

# EXHIBIT 11

**PUBLIC VERSION -  
CONFIDENTIAL MATERIAL OMITTED IN FULL**

# EXHIBIT 12

**PUBLIC VERSION -  
CONFIDENTIAL MATERIAL OMITTED**

IN THE UNITED STATES DISTRICT COURT  
FOR THE DISTRICT OF DELAWARE

RYANAIR DAC,	)	
	)	
Plaintiff,	)	
	)	
v.	)	Case No.
	)	1:20-CV-01191-WCB
BOOKING HOLDINGS INC.,	)	
BOOKING.COM B.V., KAYAK	)	
SOFTWARE CORPORATION,	)	
PRICELINE.COM LLC, and AGODA	)	
COMPANY PTE. LTD.,	)	
	)	
Defendants.	)	
_____	)	

HIGHLY CONFIDENTIAL PURSUANT TO THE PROTECTIVE ORDER

Video-recorded videoconference deposition  
of TIMOTHY J. O'NEIL-DUNNE, taken remotely on behalf of  
the Plaintiff, beginning at 9:03 a.m. and ending at  
5:24 p.m., on Tuesday, September 26, 2023, before  
JOANNA B. BROWN, Certified Shorthand Reporter No. 8570,  
RPR, CRR, RMR.

Page 2

## APPEARANCES

## FOR THE PLAINTIFF:

HOLLAND &amp; KNIGHT LLP

BY: CYNTHIA A. GIERHART, ESQ.

800 17th Street NW, Suite 1100

Washington, DC 20011

(202) 569-5416

cindy.gierhart@hklaw.com

## FOR THE DEFENDANTS:

COOLEY, LLP

BY: ALEXANDER J. KASNER, ESQ.

1299 Pennsylvania Avenue, NW, Suite 700

Washington, DC 20004

(202) 842-7800

akasner@cooley.com

COOLEY, LLP

BY: KATHLEEN R. HARTNETT, ESQ.

KRISTINE FORDERER, ESQ.

3 Embarcadero Center, 20th Floor

San Francisco, California 94111-4004

(415) 305-6527

khartnett@cooley.com

kforderer@cooley.com

## ALSO PRESENT:

THERESA MAJERS, VIDEOGRAPHER

Page 3

## INDEX

EXAMINATION OF: PAGE

TIMOTHY J. O'NEIL-DUNNE

BY MS. GIERHART 7

## EXHIBITS

PLAINTIFF'S PAGE

Exhibit 125 Expert Report of 52

Timothy James O'Neil-Dunne

(55 pages)

Exhibit 126 Booking.com B.V.'s Supplemental 52

Responses and Objections to

Plaintiff Ryanair DAC's

Interrogatory No. 2 (14 pages)

Exhibit 127 Twitter post by Mr. O'Neil-Dunne 138

(1 page)

Exhibit 128 "Advantages and Disadvantages of 201

Online Travel Agencies (OTAs)"

(16 pages)

Exhibit 129 "Web Scraping for Hospitality 226

Research: Overview,

Opportunities, and Implications"

(16 pages)

Exhibit 130 August 10, 2008, 236

"Professor Sabena's Blog"

(1 page)

Exhibit 131 August 18, 2008, 253

"Professor Sabena's Blog"

(1 page)

Page 4

## EXHIBITS

PLAINTIFF'S PAGE

Exhibit 132 April 16, 2009, 258

"Professor Sabena's Blog"

(1 page)

Exhibit 133 "Legal and Ethical Issues of 266

Collecting and Using Online

Hospitality Data" (9 pages)

Exhibit 134 "THE BUSINESS IMPACT OF 271

WEBSITE SCRAPING: IT'S PROBABLY

BIGGER THAN YOU THINK - HERE'S

WHY" (10 pages)

## UNANSWERED QUESTIONS

PAGE LINE

(None.)

Page 5

Remotely; Tuesday, September 26, 2023

9:03 a.m.

(TIMOTHY J. O'NEIL-DUNNE,

deponent, was sworn and examined

and testified as follows:)

THE VIDEOGRAPHER: We are now on the record.

This begins Media File No. 1 in the deposition of

Timothy O'Neil-Dunne. We are taking this deposition

remotely via Zoom in the matter of Ryanair DAC versus

Booking Holdings, Inc., et al., filed in the

U.S. District Court for the District of Delaware,

Case No. 1:20-CV-01191-WCB.

Today is September 26, 2023. The time is

9:03 a.m. I am Theresa Majers, the videographer, from

Magna Legal Services. The court reporter is

Joanna Brown also with Magna Legal Services.

Counsel, please introduce yourselves and who

you represent starting with the questioning attorney.

MS. GIERHART: Cynthia Gierhart on behalf of

Ryanair.

MR. KASNER: Alex Kasner of Cooley, LLP, on

behalf of defendants. I have with me Kathleen Hartnett

also of Cooley, LLP, on behalf of defendants,

Page 114

1 have a technical ability to distribute their products  
2 through GDSs for a typical impediment. There were  
3 other airlines who chose not to participate even though  
4 they could do so.

5 Q And did those airlines have content available  
6 through some computer system where a travel agent could  
7 use some technology to go access it?

8 MR. KASNER: Object to the form.

9 THE WITNESS: So pre-Internet, if you had an  
10 airline who didn't have access to the GDS and there was  
11 no Internet, there was only one -- there were only two  
12 other possible forms of access: one was to walk into a  
13 physical location, and the other one was by the phone.

14 BY MS. GIERHART:

15 Q All right. So I think, then, it's -- it's  
16 sort of a post-Internet scenario where a travel agent  
17 would use some technology to access an airline's  
18 content outside of a GDS or through an agreement with  
19 the airline direct?

20 MR. KASNER: Object to the form.

21 THE WITNESS: Okay. So there's been --  
22 there's a lot of theory about did the airlines'  
23 websites make the Internet very strong, or did the  
24 Internet make the airlines' websites very strong?

25 I think they are symbiotic. So the emergence

Page 116

1 without the airlines' permission; right?

2 MR. KASNER: Object to the form.

3 THE WITNESS: That is correct.

4 BY MS. GIERHART:

5 Q Okay. So there's not a history for decades of  
6 travel agents selling airlines' tickets without their  
7 permission?

8 MR. KASNER: Object to the form.

9 THE WITNESS: That's not actually true. There  
10 was a very robust gray market where travel agents who  
11 did have the right to issue tickets distributed theirs  
12 to subagents and had a wider network than their own  
13 agents. There's a term used which is called  
14 "consolidators" or "consolidation," and consolidators  
15 would actually sell to large numbers of travel agents.  
16 And the airlines actually liked it because it was a way  
17 for them to sell unused or wasting inventory to a  
18 broader market that wasn't necessarily fully just a  
19 regular travel agency.

20 BY MS. GIERHART:

21 Q So it was a bit of off-market selling, but the  
22 airlines kind of supported it regardless?

23 MR. KASNER: Object to the form.

24 THE WITNESS: Some. I said some. And I want  
25 to be very cautious that that -- that marketplace

Page 115

1 of the Internet as a -- as a channel of distribution  
2 and the emergence of airlines being able to distribute  
3 their content outside of the GDS sort of occurred about  
4 the same time.

5 BY MS. GIERHART:

6 Q Okay. All right. I think, just to wrap up  
7 this Section 4, if we could actually just look at this,  
8 on page 5 of your report, actually going up, there's a  
9 little summary for each section. So I just want to  
10 look at the summary for Section 4. Okay. So that's  
11 page 5, the paragraph that starts "First."

12 So you say "The development of modern airline  
13 ticket distribution -- split between direct and  
14 indirect distribution can be traced through decades of  
15 history."

16 So I think you are saying here airlines have  
17 been selling flights through travel agents and -- for  
18 decades. Is that, sort of, the first conclusion?

19 A Yes. I can wax lyrically about how that  
20 occurred; but, yes, that's for airline tickets. Yes,  
21 there were decades of history where travel agents sold  
22 tickets on behalf of the airlines.

23 Q And that doesn't distinguish between -- well,  
24 you just said selling tickets on behalf of the airline.

25 So this doesn't involve selling tickets

Page 117

1 existed. It existed far more in places like Europe and  
2 Asia than it did in the U.S., but it was definitely  
3 something that happened.

4 BY MS. GIERHART:

5 Q And your report does not make a distinction in  
6 this early period of a long history of travel agents  
7 selling tickets on behalf of airlines.

8 Your report doesn't make a distinction between  
9 authorized and unauthorized sales; right?

10 MR. KASNER: Object to the form.

11 THE WITNESS: What I was addressing in my  
12 report -- I think I was very explicit about this -- was  
13 really about the U.S. market and about the more  
14 conventional, mainstream ways that airline tickets were  
15 distributed. So I didn't want to go into that level of  
16 detail.

17 There are many, many different exceptions that  
18 are possible in these cases. I have tried to avoid  
19 doing that so that we would have a report that was only  
20 25 pages as opposed to one that could be several  
21 hundred pages long.

22 BY MS. GIERHART:

23 Q Right. Okay. But -- so the -- you are saying  
24 you were focusing on the mainstream. Sort of a  
25 standard protocol of the time was -- in the

Page 118

1 pre-Internet era was for travel agents to work  
2 through -- to sell airline tickets through approved  
3 methods, which were through GDSs, which had agreements  
4 with airlines.

5 Is that -- that's the standard procedure at  
6 that time?

7 MR. KASNER: Object to the form.

8 THE WITNESS: That was the way that it  
9 actually happened. That would be a good  
10 characterization, yes.

11 BY MS. GIERHART:

12 Q Okay. Okay. So your report, sort of, focuses  
13 on the mainstream, authorized sales of airlines through  
14 travel agents?

15 MR. KASNER: Object to the form.

16 THE WITNESS: With the caveat that the way  
17 that I described this was I focused on the  
18 United States market, which, as I said, is different  
19 from different countries in different jurisdictions.

20 BY MS. GIERHART:

21 Q Right. Okay. But that was the case in the  
22 U.S. market?

23 A The majority, yes. I mean, there were  
24 consolidations in the U.S. too, and they are fairly  
25 well known. Some of them still exist.

Page 120

1 THE WITNESS: You, kind of, concatenated lots  
2 of things. Can you unpack that maybe a little bit?

3 BY MS. GIERHART:

4 Q Sure. I think in the -- in the -- in your  
5 opinion about screen-scraping technology having been  
6 used for decades by travel agents, the early form  
7 pre-Internet was limited to manipulating content that  
8 travel agents already had from the GDS; is that right?

9 MR. KASNER: Object to the form.

10 THE WITNESS: I can answer the question this  
11 way: because the travel agents used only the GDS as  
12 their primary form of interactive system, therefore,  
13 the manipulation of screen emulation was based upon the  
14 data that was entirely controlled from the GDS down to  
15 the terminal. There was nothing between the two.

16 When you had a PC and you could emulate the  
17 screen, then the possibility started to exist because  
18 there was now an intelligence that sat at the remote  
19 end of the communication, and that was screen emulation  
20 initially, and, then, it went to screen manipulation.

21 BY MS. GIERHART:

22 Q Right. Okay. I guess I was just making the  
23 distinction, it didn't involve yet going out and  
24 getting content from somewhere else?

25 MR. KASNER: Object to the form.

Page 119

1 Q I'm sorry. A consolidator?

2 A Sorry. These -- these companies who sell  
3 tickets through different forms rather than the  
4 conventional form, there are a number of these forms  
5 that now exist for different reasons. There has always  
6 been alternative forms, and -- but when we are talking  
7 about, as I put in the report, the development of the  
8 modern airline ticket system, it was conventionally  
9 seen as being airline talked to GDS, talked to travel  
10 agent, and, then, customers came in and bought tickets  
11 from the travel agent. That was the indirect form of  
12 distribution.

13 Q Okay. Okay. And, then, this, sort of, second  
14 part of this paragraph is also, then, talking about  
15 what we were just talking about, I think. With direct  
16 distribution and indirect distribution, they have  
17 existed for at least a half a century with the latter  
18 using screen emulation and screen-scraping technology  
19 for decades.

20 So, the second point, we are talking about,  
21 sort of, the advent of screen-scraping technology.  
22 This was an early form where it was manipulation of  
23 content that the travel agents already had access to  
24 through their relationship with the GDSs; right?

25 MR. KASNER: Object to the form.

Page 121

1 THE WITNESS: I'll disagree with that because  
2 the ability to bring other forms of information, again,  
3 pre-Internet was there were physical sources of  
4 information that existed. So someone might, for  
5 example --

6 Let me try and use an example for this. I  
7 might say that, here in this travel agency, I don't  
8 want you to use United Airlines, and that was a  
9 directive from the management of the travel agency even  
10 though the travel agency was perfectly able to do that.

11 It was possible through manipulating the  
12 screen to eliminate United Airlines out of the  
13 displays. That's a possibility. That's a good  
14 scenario to illustrate how manipulation of the data  
15 that was coming to and from could actually change what  
16 was happening at the agents' workstation.

17 BY MS. GIERHART:

18 Q Right. I understand that as a removal of  
19 content. I just meant obtaining raw data. Where the  
20 content comes from was from the GDS. That's all I  
21 meant. These early, you know, stages of  
22 screen-scraping technology didn't yet involve -- and I  
23 think we said because of the example with Russia, or  
24 you only sell by phone by GDS. By going in person,  
25 there was -- there just wasn't yet a way for travel

Page 122

agents to get content before the Internet if it didn't --

A Actually, there was.

MR. KASNER: Object to the form.

THE WITNESS: Sorry.

Actually, there was. There's a service on which the GDSs are based, which is called the "OAG." The OAG stands for the "Official Airlines Guide." And the OAG publishes schedules for most airlines whether they participate or not. Therefore, I could have access -- and I, in fact, remember doing this -- where you could bring in the information on dial-up to the OAG source, and you could see that there was, in fact, other airlines who were operating at that particular moment on a particular route, who might not be participating inside a GDS.

And, therefore, you would now be able to say "Ah, I see that that can happen. How do I get access to that" -- "How can I make a booking?"

Well, it will say -- you won't know that from the OAG because the OAG only published schedule data. It didn't publish any inventory data.

And, therefore, I can go and say "Ah. But I can see how this flight exists. Let me see how I can get access to it, probably by the phone."

Page 124

Q Right. Right. Okay. All right. I'll start going just a little bit into Section 5, and, then, we can take a break. Okay.

So Section 5 is now -- so, now, we are moving into the "MODERN METHODS FOR BOOKING CONSUMER TRAVEL." It starts on page 9 of your report.

And, yeah, actually, just to clarify because, when I was reading this, I just wanted to make sure, this first paragraph under -- under Section 5 at the very top of the page, it says -- this section begins by discussing, sort of, direct and indirect distribution, which I see, in Section 5, it says -- it then discusses how an airline's business model will drive an airline's decision to encourage and pursue a particular mode of distribution. I think that's Section 6.

Can you confirm if that's what you meant?

A Let me have a quick look. Because it is a comparison to -- when I get into Section 6, I go more into the details, which is the airline distribution preference. And, then, I break that out into the full-service and low-cost carriers. So I try to focus Section 5 on probably the context of indirect distribution and then into the detail because airlines were more important for the context of this report is why I created the whole extra section in 6.

Page 123

MS. GIERHART: Okay.

THE WITNESS: In fact, there are ways to manipulate the GDS to actually show what is known as schedule data only -- sorry -- schedule data only, which is possible if you know how to manipulate the GDS.

BY MS. GIERHART:

Q So it would be the equivalent, say, of getting flight itinerary information or route information but not -- there was no ability, you know, to book flights or issue tickets. Then, you said, you'd have to call or, you know --

A Yes.

Q -- use some other method?

MR. KASNER: Object to the form.

THE WITNESS: That's a good way to describe it.

MS. GIERHART: Okay.

THE WITNESS: But I just would caution you about the, no, not possible because there was always somebody somewhere in some garage who was tinkering and trying to do things. So I would just say it was definitely not mainstream. Mainstream was the GDS activity.

BY MS. GIERHART:

Page 125

Q Okay. So in this -- in this -- all right. So 5.1, we are talking about direct distribution, and you note in here -- I'm trying to find the exact wording. All right. The second sentence, "This channel gives airlines the most control over their pricing and marketing, and it allows them to collect customer data that can be used to improve their loyalty programs and targeted marketing campaigns."

Why -- we talked about this a little before, but why might an airline want to control their pricing and marketing?

MR. KASNER: Object to the form.

THE WITNESS: For many different reasons. Targeting is, as I've put in here, targeting of products and services and marketing.

The loyalty issue is a big one because, prior to the Internet, loyalty was not as precise a science as it is today. So the target -- the control is something that the airlines wish to exert because, if they feel that they can keep a captive customer, they feel that they can leverage that customer to become more valuable to them as a supplier.

BY MS. GIERHART:

Q Why might an airline want to collect customer data?



Page 218

1 talking about OTAs being beneficial for the purpose of  
 2 comparing flight information, not necessarily booking?  
 3 But -- so when you are talking about  
 4 information-sharing, aggregation, is that all just  
 5 related to comparing flights?

6 MR. KASNER: Object to the form.

7 THE WITNESS: In general, yes, but it also  
 8 applies to other forms of travel information like  
 9 hotels.

10 BY MS. GIERHART:

11 Q Okay. But it's not referring to a benefit to  
 12 the customer of booking flights through an OTA; right?

13 A This section is just dealing with the  
 14 information-sharing and aggregation, those two things  
 15 brought together.

16 Q Okay. So -- okay. So you -- I just want to  
 17 clarify and make sure. So the opinion does not  
 18 indicate any reason why booking a flight on an OTA is  
 19 preferable to booking a flight through an airline?

20 A If you look at the penultimate paragraph under  
 21 7.2 on page 15, "From a consumer (i.e., traveler)  
 22 perspective, OTAs can provide a powerful  
 23 information-gathering and presentation tool, helping  
 24 consumers by compiling and comparing different flight  
 25 itinerary and fare information," also known as offers.

Page 220

1 out -- that's why I prefer to use that term -- which is  
 2 the process of -- it's automating the process of  
 3 accessing a web service and recording or  
 4 scraping/pausing the information displayed on the  
 5 screen.

6 So, when we are looking at the website, I'm  
 7 looking at the data as opposed to all of the different  
 8 pixels.

9 Q Okay. It seems limited to visual data that  
 10 you'd see on the website, but, then, you scrape or  
 11 pull; is that correct?

12 A Yeah. It's --

13 MR. KASNER: Object to the form.

14 THE WITNESS: Sorry. It's the data as opposed  
 15 to whether it's green screen, green type, white on  
 16 green, or yellow on black. Yes. It's related to the  
 17 data that's in there rather than just the pixels.

18 BY MS. GIERHART:

19 Q Okay. And in relation to OTAs and airlines  
 20 referring to screen scraping here, present day, would  
 21 refer to pulling light itineraries and routes, maybe  
 22 ancillary options, pulling all of that data?

23 That's all screen scraping; right?

24 MR. KASNER: Object to the form.

25 THE WITNESS: So screen scraping this, and as

Page 219

1 So that isn't talking about booking. This  
 2 particular section does not refer to booking. That is  
 3 correct.

4 Q Okay. Okay. And, then, going to the next  
 5 section regarding screen scraping as an information  
 6 aggregation tool, so now we get to the question of how  
 7 do you define "screen scraping" as used in this report?

8 A Right. So I tried to come up with a succinct  
 9 element here. So screen scraping -- I actually prefer  
 10 the term "screen pausing" because, if you think about a  
 11 screen, going back to our emulation story, when we were  
 12 dealing with emulation, we were only dealing with  
 13 characters/text that was on the screen.

14 When you move to a PC in today's highly dense  
 15 pixel 4K TV screen and 4K computer screen, 104- --  
 16 1080P or whatever this is, you could be looking at the  
 17 entire screen and say "Well, am I looking at just the  
 18 screen with every pixel, every color, or whatever; or  
 19 am I looking at the information?"

20 So screen scraping could encompass the whole  
 21 of that screen including the bits down at the bottom.  
 22 I just see it on my screen. I've got a weather app.  
 23 I'm not going to use that because that's not really  
 24 relevant to us.

25 Screen pausing is taking data that is relevant

Page 221

1 I put in here, the process extracts data from a  
 2 computer screen by reading the pixels, starting at that  
 3 level, on the screen and converting those into the --  
 4 and getting rid of the unnecessary stuff into text or  
 5 other data formats that can then be used.

6 So, in this particular case, you described a  
 7 series of elements that one can find on an airline  
 8 website such as flight information, times, or whatever  
 9 it might be and pulling that into and then making that  
 10 available.

11 BY MS. GIERHART:

12 Q Okay. And would you call it something  
 13 different, the process of -- of OTAs actually booking  
 14 the flights?

15 Connecting to the API and booking a flight, is  
 16 that called -- that's not screen scraping, but that's  
 17 something else?

18 That's booking a flight; is that accurate?

19 MR. KASNER: Object to the form.

20 THE WITNESS: So consider that you -- we are  
 21 doing two different things here. So the first thing,  
 22 screen scraping, per se, is just grabbing the data,  
 23 bringing it into an area that makes it usable. That's  
 24 screen scraping.

25 When you are doing something else like you are



Page 222

manipulating the data or doing something with it as I describe in some of the elements below, then that's different. Then I'm now processing the information in the same way that I would if I was a human on the machine.

BY MS. GIERHART:

Q Well, there's processing information, and, then, there's also -- or there's also connecting to an API to book a flight.

Would you consider that part of the screen-scraping process or something else separate?

MR. KASNER: Objection to the form. Testifying. Foundation.

THE WITNESS: So I think what I've described is you need to get a hold of information before you can do something with it. Once I've got the information into a work -- let's call it a work area -- then I can do something with that data.

So, if I'm going to use it and redisplay it, that's an action. Then the screen-scraping effort is done. I haven't manipulated the data, or I might have manipulated the data and put it somewhere, but I'm not interacting with the place where I got the data.

MS. GIERHART: Okay. I think I understand.

THE WITNESS: See, an API -- when you describe

Page 224

Q Uh-huh.

A -- grabbing the information from different sources using screen scraping as one of the techniques, that's what I'm trying to get across.

Q Okay. You go through the key benefits.

Did you consider the downsides of scraping in forming your opinion?

A Yes.

MR. KASNER: Objection to the form. Foundation.

BY MS. GIERHART:

Q Okay. What downsides did you consider?

A There's many different downsides -- two potential downsides. Whether or not those downsides are particularly important or not depends on the use. So I think it sort of depends. The question, I think, overall is to consider the benefits in total as aggregate -- sorry -- not aggregate but unbalanced.

Is screen scraping a useful capability to bring information in that can then be manipulated by the human or by a machine?

Screen scraping is a very common tool for doing that, and I described that.

But, yes, I did consider alternatives, which is not just there are downsides to it, which can be

Page 223

an API, an API is essentially -- the key term is the "I." That's the interface. So an application programming interface means it's bidirectional. So I can get stuff, and I can put stuff.

When I'm screen scraping, screen pausing, I'm grabbing the data. I'm only doing one of those functions, whereas, an API, you have to do both.

BY MS. GIERHART:

Q Right. Okay. And so I think -- I just wanted that understanding going into, then, the next part of the report. For example, it talks about screen scraping offers several key benefits. So I just wanted to understand that when you are saying screen scraping offers benefits, we are talking about going and pulling data from airlines and displaying -- well, maybe not even displaying it. That might be a Step 2.

But the act of pulling data from an airline's website is what you are talking about when you -- when you say "Screen scraping offers several key benefits"; is that right?

A That's correct.

Q Okay.

A In fact, if I can, sort of, look at the title of 7.2.1, which is screen scraping as an information aggregation tool --

Page 225

varied.

Q Okay. I think while -- I'm going to go to Exhibit -- the next exhibit, which is your article referenced in Footnote 14, the second one, the Han and Anderson article --

A Uh-huh.

Q -- which is going to be Exhibit 129.

(Deposition Exhibit 129 was electronically received and marked for identification.)

BY MS. GIERHART:

Q Okay. Do you see that?

A I'm there. Where would you like me to go?

Q Okay. Can you go to page 102 of the article.

A Well, this was -- hold on a second because the way it's coming up is I've only got to page 16.

Q Oh, sorry. It's page 14 of 16.

A Okay.

Q It says 102 at the top.

A Oh, I see. Yes. Understood.

Q Okay. So this page of the article talks about "Ethical Concerns."

Do you see that?

A I do.

Q Actually, sorry to make you jump. I should just go to the title of the article so we know -- have

Page 226

1 some context. So the article is entitled "Web Scraping  
2 for Hospitality Research: Overview, Opportunities, and  
3 Implications." So, in that context, we are talking  
4 about web scraping and ethical concerns related to it.

5 Do you remember reading this when you were  
6 writing your report?

7 A Yes. And -- I do.

8 Q Okay. Do you see, at the bottom of the first  
9 paragraph there in that section, it says: "However,  
10 web scraping information that is accessible exclusively  
11 to the members and requires logging in is illegal, as  
12 this behavior explicitly violates the terms of  
13 service"?

14 A Hold on a second. I'm sorry. I can't find  
15 that. I remember the issue. I just can't see it.  
16 Where is it? We are on the same page?

17 Q It's in the middle of the page on the  
18 left-hand --

19 A Oh, yeah, I see it. So -- yes. "However, web  
20 scraping information that is accessible exclusively,"  
21 yes, I see it there.

22 Q Do you agree with this?

23 A So I went to look at the original  
24 Supreme Court ruling. My interpretation of it -- see,  
25 I'm not a lawyer -- said that, in general, it's --

Page 228

1 which was whether or not it's required. So I looked at  
2 actually going in to Ryanair and seeing what Ryanair  
3 did on its website. And, in fact, I show -- at the  
4 back of my report, there are some sample screens that  
5 show that I went in. And, with my assistant, we  
6 actually looked at two scenarios where we went in using  
7 Ryanair direct and where we went in using Booking.com  
8 to see if there was really a difference between the  
9 two.

10 Q A difference in what?

11 A In how we got to the data and whether the  
12 customer was impacted at the end of that.

13 Q Okay. And when you refer to getting the data,  
14 so we are talking about anything leading up to booking  
15 the ticket but not including booking the ticket?

16 MR. KASNER: Object to the form.

17 THE WITNESS: That is correct.

18 BY MS. GIERHART:

19 Q Okay. Actually, did your report take an  
20 opinion at all regarding anything at the moment of  
21 booking a ticket or after booking the ticket?

22 MR. KASNER: Object to the form.

23 THE WITNESS: I think we described what  
24 happens in the screenshots to show that we actually  
25 went through that, and we talk about myRyanair as part

Page 227

1 there is legalized screen scraping; right?

2 So the question of this particular question --  
3 this particular issue about whether this behavior  
4 violates the terms of service or not, I think, is an  
5 interpretation that these guys who wrote the article  
6 had. And since web scraping -- screen scraping is such  
7 a common behavior, if this was the case, I think there  
8 would be a lot more cases where this behavior was  
9 deemed to be illegal because this is not the only place  
10 that does it. Web scraping occurs in many other  
11 industries, and I think I cite other articles where it  
12 takes place. Travel just happens to be one of them  
13 because --

14 Q I -- I was just going to say I think the  
15 distinction here is just mentioning scraping where  
16 there's a login and violating the terms of service.

17 So I'm just wondering if that's a  
18 determination that you considered or not in writing  
19 your report.

20 MR. KASNER: Objection to the form and  
21 testifying.

22 THE WITNESS: Yes, I did.

23 BY MS. GIERHART:

24 Q Okay. What did you consider about that?

25 A I considered several things in here, one of

Page 229

1 of the -- it's in the report.

2 BY MS. GIERHART:

3 Q Okay. I guess I didn't phrase it very well.  
4 I guess, just generally, your report, when it talks  
5 about benefits of screen scraping or aggregating data,  
6 OTAs being able to show different flights, that's  
7 really all talking about the search process leading up  
8 to but not including booking a ticket; is that right?

9 MR. KASNER: Object to the form.

10 THE WITNESS: Most of the piece is concerned  
11 with web scraping, which is grabbing information.  
12 Remember, that's a one-way kind of way to look at it.  
13 That's pulling the information in. That's not  
14 interacting with it out.

15 BY MS. GIERHART:

16 Q Okay. Okay. All right. So I just want to  
17 look at another -- on the same page, but if you move to  
18 the right column, it's about five lines down. The  
19 author is saying: "A common ethical concern regarding  
20 web scraping is related to the problem of sending too  
21 many requests to the host over a short span of time"  
22 and that [as read] "A typical web scraper involves  
23 querying a website repeatedly. If overused, the  
24 practice can prevent others from accessing the  
25 website."

Page 230

1 Are you familiar with this practice?

2 A I am very familiar with this practice.

3 Q Okay. And what -- how are you familiar with  
4 it?

5 A The behavior of lots of requests coming in is  
6 a common challenge that a website has to deal with.  
7 So, for whatever reason, floods of requests come into a  
8 website, and they can be driven by a wide range of  
9 different reasons. It could be I'm running a  
10 promotion, which could stimulate traffic. It could be  
11 that there is something known as a DDOS, a  
12 denial-of-service attack. It could be that somebody is  
13 trying to scrape my site if I'm the site itself.

14 So having been on the side of someone who has  
15 been the target of scraping, you build tools to support  
16 and ensure that you don't have it, and you hope that  
17 whoever is on the other end, who is intending to scrape  
18 you, is behaving ethically -- sorry. Ethically is a  
19 judgment decision -- is behaving well so that you don't  
20 actually do things like that to prevent a customer from  
21 coming in. I'm familiar with this practice.

22 Q Is it a known, I guess, harm to airlines'  
23 websites --

24 Is web scraping known to harm airlines'  
25 websites?

Page 232

1 A Yes. And my experience definitely points me  
2 in that direction. I've been -- I've been at this -- a  
3 victim, if you like, of a website I was responsible  
4 for, receiving lots and lots of requests from both  
5 scrapers and from bad actors.

6 Q Which air -- which website was that?

7 A I believe I alluded to this earlier when we  
8 were running the earlier lines, and we had a portion of  
9 the website on each one of seven airlines. So our  
10 particular section was the one getting hammered.

11 Q Sorry. Was this -- I missed it.

12 Was this Flair Airlines or with a different --

13 A No. And one of the things in my CV, it talks  
14 about a company called Air -- excuse me -- Air Black  
15 Box, and Air Black Box was the technology provider to a  
16 series of airlines. There were seven of them in all.  
17 And we ran a portion of the website that created a  
18 particular function, and we ran that function for each  
19 one of the airlines. So we didn't just have to deal  
20 with one airline's website. We actually dealt with  
21 seven and an eighth, which was our own, which was  
22 ValueAirlines.com. So we were very familiar with these  
23 types of behavior.

24 Q Is there a reason why you didn't include these  
25 disadvantages in your report and weigh them and then

Page 231

1 MR. KASNER: Object to the form.

2 THE WITNESS: It's known to be capable to harm  
3 airlines' websites. I chuckle at this because, in the  
4 early days of the Web, it was a real problem. Now we  
5 have tools which identify this type of request, and if  
6 you are a well-managed website, you can figure out ways  
7 to eliminate that so that you don't have to deal with  
8 it whether these unnecessary requests are coming in.  
9 So you can monitor it. You can determine it. Those  
10 tools exist and are well deployed by every competent  
11 website owner.

12 BY MS. GIERHART:

13 Q Okay. Would you agree there's some expense in  
14 developing those tools and constantly fighting scraping  
15 activities?

16 MR. KASNER: Object to the form.

17 THE WITNESS: I have more stories I can tell  
18 you about having to do it, and most, if not all,  
19 website owners, whether it's airline or otherwise,  
20 should expect this behavior. And most of the website  
21 tools now give you a way to handle it, and it's  
22 cost-effective.

23 BY MS. GIERHART:

24 Q And did you consider these downsides when  
25 forming your opinion in this report?

Page 233

1 ultimately conclude that, you know, one outweighs the  
2 other? because they are just absent from the report.

3 MR. KASNER: Object to the form. Testifying.

4 THE WITNESS: I felt that the best way to talk  
5 about this was to have other people describe it. I've  
6 got -- as I said, I've got personal experience of doing  
7 this. From this experience and my understanding of how  
8 airlines work, it's standard practice that you expect  
9 this type of behavior to occur. So I didn't think it  
10 was necessary to go beyond that because you have an  
11 article such as this which describes those types of  
12 conditions.

13 BY MS. GIERHART:

14 Q Right. I mean, you just have to do a lot of  
15 digging to go find it, and there's no indication in  
16 your report that you are going to --

17 You didn't say "See disadvantages in  
18 Footnote 14." It just says "screen scraping is a  
19 process of accessing a web service. See Footnote 14."

20 MR. KASNER: Objection to the form.  
21 Argumentative.

22 THE WITNESS: I think the purpose of quoting  
23 this article as I've tried to do, in other words, is  
24 there is a lot of information available to explain to a  
25 user -- I'm sorry -- a reader that these issues do



Page 234

1 exist. I hoped that by quoting esteemed journals like  
2 this particular one of the value that that brought.  
3 Cornell Hospitality School is very well known, and this  
4 publication is pretty authoritative; right?

5 I don't have to agree with everything, but  
6 it's a good layout of exactly what happens. I mean,  
7 they are talking mostly about hospitality sites, but it  
8 does apply to airline sites. And I felt this was a  
9 good way to illustrate the challenge by quoting  
10 something from an authoritative journal.

11 MR. KASNER: Cindy, at the end of this line of  
12 questioning, can we take a break?

13 MS. GIERHART: I think we can do a quick break  
14 now and then come back. Another ten minutes?

15 MR. KASNER: Okay.

16 THE WITNESS: Okay.

17 THE VIDEOGRAPHER: We are off the record at  
18 3:52 p.m.

19 (Off the record.)

20 THE VIDEOGRAPHER: This is the beginning of  
21 Media File No. 6. We are back on the record at  
22 4:02 p.m.

23 BY MS. GIERHART:

24 Q Okay. So we were just talking about -- in the  
25 report, you mention the benefits of screen scraping,

Page 236

1 still seems relevant surprisingly to this day.

2 MR. KASNER: Objection to the form.

3 Testifying.

4 BY MS. GIERHART:

5 Q Do you remember this article -- or this blog  
6 post?

7 A I remember writing it. I don't remember a lot  
8 of the detail. So there's a lot in here.

9 Q Okay. Well, first, it starts off saying  
10 Ryanair is going after screen scrapers. So the second  
11 paragraph there says "So it is not surprising that  
12 occasionally they go after intermediaries who are  
13 re-selling the Ryanair product."

14 So I guess, first, the question is just, it  
15 was known in the industry since at least 2008 that  
16 Ryanair didn't want screen scrapers on their website;  
17 is that right?

18 MR. KASNER: Object to the form.

19 And if you need a second to read the article,  
20 please do.

21 THE WITNESS: Yeah. Let me read the whole  
22 article.

23 MS. GIERHART: Yeah. Just let me know when  
24 you are done.

25 THE WITNESS: Yeah. Yeah, I remember writing

Page 235

1 and you indicated that you also considered the  
2 disadvantages of screen scraping. I want to come back  
3 to that for a second, if we can switch gears, and I  
4 wanted to look at some of your blog posts of yours.

5 So I can share -- this is going to be

6 Exhibit 130. You might have to zoom in a little bit.

7 (Deposition Exhibit 130 was electronically  
8 received and marked for identification.)

9 BY MS. GIERHART:

10 Q First, just out of curiosity, I notice your  
11 blog is called -- and your Twitter handle is  
12 Professor Sabena. Where did that come from?

13 A So I tend to have a somewhat professorial way  
14 that I speak to people in situations. So my nickname  
15 was "Professor." And when I was trying to come up with  
16 a handle, I wanted to have the antithesis of something  
17 that might appeal to people who were aviation, in other  
18 words, AV geeks and people who knew about the airlines,  
19 and Sabena was probably the world's worst airline. In  
20 fact, its acronym stands for "such a bloody experience,  
21 never again." So Professor Sabena is my nickname, my  
22 nom de plume, as it were. This -- I see this is from  
23 what I wrote in 2008.

24 Q Yes, it is from a while ago, but it is talking  
25 about "Ryanair clamps down on Screen Scrapers." So it

Page 237

1 this. I'm finished.

2 BY MS. GIERHART:

3 Q Okay. So, to start, it seemed, in the  
4 industry, it was known that Ryanair did not approve of  
5 screen scrapers; is that right?

6 MR. KASNER: Objection to the form.

7 THE WITNESS: Sorry. It's been well known for  
8 a long time that Ryanair does not like screen scrapers,  
9 and it is well known that other airlines don't like  
10 screen scraping.

11 BY MS. GIERHART:

12 Q Okay. And why would they not like screen  
13 scraping?

14 A So in this article, which, remember, is very  
15 old -- it is 15 years old. That's a lifetime, maybe  
16 generations inside the travel industry, inside the Web  
17 for sure. So, at the time, this was something that  
18 was, I think, more of a challenge. It was still a  
19 fairly new phenomenon at the time about screen  
20 scraping. So screen scraping was not -- was fairly  
21 novel at this time. I think that's the best way you  
22 can describe it.

23 Q Okay. And as far as why -- why would airlines  
24 at this time not like screen scraping?

25 A I actually put in here that there was actually

Page 238

1 some risk. Where is the bit in here? I read that. It  
2 says -- if you go to the paragraph, it says "I have an  
3 issue with the use of Screen Scraping" --

4 Q Uh-huh.

5 A Right? -- "and have never been a huge fan  
6 despite the value that can be derived"; right?

7 [As read] "The core root of the problem is  
8 that the backend systems are really not designed to  
9 cope with these types of searches and therefore it  
10 creates a resource issue which can ultimately cost  
11 money. It is fair therefore for a scraper to" --  
12 sorry. "Is it fair therefore for a scraper to make  
13 money out of somebody else's resources?"

14 So I asked that question, a rhetorical  
15 question, and I don't answer it. I just pose it as a  
16 question.

17 What is interesting is, at the time, any form  
18 of screen scraping -- remember, this is long before big  
19 DDS attacks had become prevalent. So it was just a  
20 resource issue because there wasn't the tools of  
21 sophistication.

22 Now, I think that it is completely different,  
23 and this is normal behavior to expect, screen scrapers  
24 or something particularly in travel sites to come and  
25 hit you. So the tools have been developed. So I don't

Page 240

1 inside the industry that you expect that people will  
2 use some form of screen-scraping tool, and therefore, I  
3 think it's kind of a nonissue these days. Whereas  
4 before, when I wrote this article 15 years ago, it  
5 absolutely was a big issue.

6 BY MS. GIERHART:

7 Q You were starting to say -- you were comparing  
8 from when you had to deal with it with conversations  
9 you are having with people today.

10 When did you have to deal with it?

11 A Oh, sorry. This was in the 2017, '18, and  
12 '19.

13 Q And it was a big problem then?

14 A No. It was -- it was conquered, if you like,  
15 from that time. So seven years after -- 2017. Seven  
16 years after this article, by then, I had started  
17 Air Black Box. This was from 2010. I started  
18 Air Black Box in 2012, and we had to deal with it. And  
19 it was one of the first questions I asked of my web  
20 team.

21 And I said, "Okay, how do we defend against  
22 this? How do we defend against massive amounts of data  
23 requests? How do we synthesize those data requests?  
24 How do we stop bad actors from entering?"

25 And those, by then, even in 2012, have become

Page 239

1 think that statement is true anymore.

2 Q So it's your belief that screen scraping no  
3 longer causes harm to an airline's website?

4 MR. KASNER: Object to the form.

5 THE WITNESS: I'm not saying that. I'm saying  
6 that it is not as much of an issue as it was when I  
7 wrote this article.

8 BY MS. GIERHART:

9 Q Okay. And what is that based on?

10 Have you lately worked with airlines to combat  
11 screen scraping?

12 MR. KASNER: Object to the form.

13 THE WITNESS: I would characterize it in a  
14 slightly different way, is that when I last had to deal  
15 with it in my own personal experience and in my  
16 discussions with airlines who now have to deal with the  
17 problem on their websites and other travel websites  
18 because it's not just airlines who get screen-scraped,  
19 it is also -- sorry -- OTAs who get it. In fact, I  
20 talk there about one website called "Trip.com."

21 At the time, Trip.com was owned by somebody  
22 else, and they were a bad actor. But, now, it's owned  
23 by Ctrip. It's become a good actor in that they're  
24 using Trip.com in a good way.

25 So I think it's part of the standard behavior

Page 241

1 a lot more normal because this issue was raised at this  
2 time and the people had started to deal with it. By  
3 the time we were in public in 2016 with Air Black Box  
4 and its customers, it was an absolute normal situation,  
5 and we were dealing with this as we were dealing with  
6 other things like DDOS attacks.

7 Q And do you think you were getting  
8 screen-scraped at the same level that airlines like  
9 Ryanair are?

10 MR. KASNER: Object to the form.

11 THE WITNESS: Can you give me perhaps a better  
12 definition of "level."

13 BY MS. GIERHART:

14 Q Do you think that the volume of screen  
15 scraping that your company received would be equivalent  
16 to the volume of screen scraping that Ryanair receives?

17 MR. KASNER: Object to the form.

18 THE WITNESS: Ryanair was a massively much  
19 larger airline. It's one of the top airlines on the  
20 Web anyway. So we were nothing like that. We were  
21 confined to a very specific area in Asia Pacific. So,  
22 the total volume, I would say no; but in comparison to  
23 what I understand Ryanair is getting, I don't have any  
24 data on that subject. I tried to find some, but I  
25 can't find any. I would say, probably, it was far more

Page 242

1 impactful on us than it was on Ryanair.

2 BY MS. GIERHART:

3 Q Even though it would be a much smaller volume?

4 A Yes. Because let's say we got a thousand  
5 bookings a day and we had 10,000 hits that came from  
6 screen scrapers, that would translate into -- Ryanair  
7 into literally millions of hits in comparison.

8 Q And so --

9 A We got a lot, and in one particular case, I  
10 can recall most of the traffic coming in was from  
11 people attempting to scrape, and they were actually a  
12 malicious actor in the specific case that I'm talking  
13 about because they were trying to create a copy of our  
14 website so that they could create a fake website.

15 So the detailed knowledge was not just that  
16 the actor was doing something bad if we count screen  
17 scraping as bad; but they had a malicious intention of  
18 doing that not to the benefit of our customers but  
19 diverting customers to their website so they could  
20 steal from consumers. That's what --

21 Q Okay.

22 A -- (inaudible) OTAs do.

23 Q But -- so you don't know the level of attacks  
24 that Ryanair has to deal with in the last five years or  
25 even today; right?

Page 244

1 document. I recall that it actually talks about what  
2 Ryanair provisions as tools, and it has a cost factor  
3 associated with it. And, to me, in my professional  
4 opinion, the size of the business that Ryanair does and  
5 the number of bookings and web searches that they  
6 receive, that seemed to be quite normal.

7 BY MS. GIERHART:

8 Q Is it a widely held belief in the industry  
9 that screen scraping just isn't a problem anymore?

10 MR. KASNER: Object to the form.

11 THE WITNESS: Thinking about that, obviously,  
12 I can't speak for everybody nor the industry as a  
13 whole, but I would say that screen scraping is normal  
14 behavior inside the travel industry, and people accept  
15 it and handle it.

16 BY MS. GIERHART:

17 Q Right. But none of that says that it isn't a  
18 problem. Do you agree that there is a difference?

19 MR. KASNER: Object to the form.

20 Argumentative.

21 THE WITNESS: I think, if we were to ask a  
22 slightly different question, "Do people accept that  
23 DDOS attacks are normal?" we would say the industry  
24 expects it. So we deal with it. I think the same  
25 applies to screen scraping. Although they are very

Page 243

1 MR. KASNER: Object to the form. Foundation.

2 THE WITNESS: I do not have specific knowledge  
3 of that, no.

4 BY MS. GIERHART:

5 Q Okay. So you can't really know if the problem  
6 is solved for Ryanair, can you?

7 MR. KASNER: Object to the form.

8 THE WITNESS: I believe there is a document --  
9 I don't recall which one -- which talks about what  
10 Ryanair does to combat the issue. I think they even  
11 quote it in one of the responses. They quoted the  
12 amount of money that they spend on anti-DDOS and screen  
13 scraping. It didn't appear to be a lot of money, which  
14 is what I would expect for an organization the size of  
15 Ryanair.

16 BY MS. GIERHART:

17 Q I would be curious which document you are  
18 referring to.

19 MR. KASNER: Object to the form.

20 THE WITNESS: God, I'd have to go back and  
21 look, but it was one of the responses. Do I have to --  
22 can I find it now, or should I just wait?

23 MS. GIERHART: I'd rather we move on if you  
24 don't know offhand.

25 THE WITNESS: No. I mean, I recall the

Page 245

1 different reasons, I think we just accept it.

2 BY MS. GIERHART:

3 Q I guess the question is you had the opinion in  
4 2008 that you had an issue with the use of screen  
5 scraping and you were not a fan.

6 Are you a fan today of screen scraping?

7 MR. KASNER: Object to the form. Foundation.

8 THE WITNESS: I will never say that I am a fan  
9 of screen scraping. That's not what I'm saying here.  
10 What I'm saying here -- and I think I still hold the  
11 belief -- is if there is a better way to do things,  
12 then we should go for it. Screen scraping should not  
13 be something that should be a be -- a goal in and of  
14 its own self; but if it's the only way to achieve the  
15 goal, then you just accept it, and you move on. You do  
16 what you can.

17 BY MS. GIERHART:

18 Q I wanted to go up just a little bit in the  
19 article. You talk about there are three main types of  
20 screen scraping. Type 1, you say, is aggregate and  
21 then forward the user to the supplier. So you give the  
22 example of Kayak. So I think that's your example you  
23 were talking about where you aggregate data and then  
24 link to the supplier.

25 And, then, Type 2, you say, is "screen scrape

Page 246

1 aggregate and then make the booking in the background  
2 without redirecting the consumer." Would --

3 Is Booking.com in the Type 2 category?

4 A So these types were defined by me in 2008.  
5 Since I haven't seen this article in at least 15 years,  
6 I'd have to think about that more before I could give a  
7 considered opinion, and since you are presenting this  
8 to me, I would ask for time to think about that  
9 question.

10 Q Sure. You can think about it.

11 And your response could be, you know, "In  
12 2008, this is how I would have considered it in 2008."

13 A In 2008, I think it was clearer because there  
14 was really -- there weren't as many options or as many  
15 nuances to it. So judging by the three criteria, if  
16 somebody else had written this, if I was using that as  
17 the criteria, I would say kind of, yes, that they are a  
18 Type 2, not completely yes as Type 2 because I would  
19 probably have had -- if I was looking at this now, the  
20 different types of screen scraping, I would have many  
21 more differentiations.

22 Q Okay. And I guess is the same --

23 Would Priceline fall into that Type 2?

24 MR. KASNER: Object to the form.

25 THE WITNESS: I would answer it the same way.

Page 248

1 BY MS. GIERHART:

2 Q Okay. I guess, going into your thinking  
3 there, Type 2, two paragraphs down, you say "Type 2 is  
4 the nasty one because it is essentially fooling the  
5 customer."

6 You do mention Ryanair here. Then you also  
7 say "It is also EXPRESSLY verboten with Ryanair and  
8 just about everyone else." And, then, you go into the  
9 terms and conditions for Ryanair.

10 A Correct.

11 Q So it sounds -- was your belief in 2008 that  
12 this Type 2 where the screen scraper is making the  
13 booking on the website without redirecting, that those  
14 were, sort of, this -- I mean, you clearly had strong  
15 feelings against them; is that accurate?

16 MR. KASNER: Object to the form. Misstates  
17 the document.

18 THE WITNESS: I'd like to agree with counsel.  
19 I think you are misstating the point here. What I am  
20 trying to say here is that -- that for the Type 2 -- I  
21 tried to focus very clearly on Type 2, and I was  
22 expressly thinking about these two players, V-tours,  
23 who I don't think exists anymore, and Bravofly that  
24 was, kind of, a cross-seller [ph].

25 And Bravofly was well known for basically

Page 247

1 MS. GIERHART: Okay.

2 THE WITNESS: And if you look there, I say  
3 that -- I name examples here of three -- of two  
4 players, V-Tours and Bravofly. And at the time -- and  
5 my memory is going to be a bit sketchy here because I  
6 haven't thought about this for a while, but Bravofly,  
7 if I recall correctly, was a well-known bad actor of  
8 Type 2, which is why I called them out because they did  
9 not provide any indication to the customer as to what  
10 they were doing and they had lots of problems and it  
11 was, kind of, well known.

12 The issue, I think, is the whole industry has  
13 moved on; and therefore, it doesn't appear to be the  
14 same as the situation that visited in 2008. I think it  
15 would be unfair to label either Booking or Priceline or  
16 even any other player you care to name in the same  
17 category without considering them in much more detail.

18 So I would just say I tend to think that there  
19 were very specific players in 2008, and I named two of  
20 them who I think were bad actors at the time. Whether  
21 or not they cleaned up their act, I can't say.

22 Does -- do I think that Booking and  
23 Priceline -- Booking.com and Priceline exactly fall  
24 into this category? I'd have to think about that.

25 ///

Page 249

1 hoodwinking the customer, not just in the area of  
2 screen scraping but in other areas, and they acquired a  
3 very bad reputation for doing that inside the industry.

4 So my first categorization is, yes, if you are  
5 completely fooling the customer, which I believe that  
6 Ryanair -- I'm sorry. Excuse me. Not Ryanair. Excuse  
7 me -- that Bravofly was. And I don't think the same  
8 thing applies, which is why I think that Booking.com  
9 and Priceline don't fit the mold for the strict  
10 interpretation I had at the time with Type 2.

11 Again, this is reinterpreting something that I  
12 had from 15 years ago.

13 BY MS. GIERHART:

14 Q Okay. I guess, to scroll down to the  
15 paragraph we were looking at before where it starts off  
16 "I have an issue with" and it ends with that you said  
17 there's a rhetorical question, is it fair, therefore,  
18 for a scraper to make money out of someone else's  
19 resources?

20 Do you have an opinion on that?

21 MR. KASNER: Object to the form.

22 THE WITNESS: I think that's still valid.

23 That's still a valid statement because, again, going to  
24 the specific case that we are talking about, Bravofly,  
25 which is the one I remember, Bravofly was absolutely



Page 250

1 making money and was doing -- I think Michael O'Leary  
2 calls them "pirates," the pirate behavior of charging a  
3 completely different price to the one that Ryanair was  
4 putting out.

5 So the issue was different, and I think,  
6 therefore, that the statement is true that, in this  
7 particular case, the scraper was making money by  
8 misrepresenting things to the consumer and using the  
9 resources of the airline because they were charging a  
10 significant premium.

11 BY MS. GIERHART:

12 Q Okay. So if an OTA substantially marked up an  
13 airline's flights and was able to sell that flight  
14 through screen scraping through using the airline's  
15 resources, that would not be fair in your opinion?

16 MR. KASNER: Object -- object to the form.  
17 Calls for speculation.

18 Can I answer?

19 MR. KASNER: Yes.

20 THE WITNESS: So I would add one component,  
21 and was not being clear, as was the case with Bravofly,  
22 purporting to be something that they were not.

23 BY MS. GIERHART:

24 Q Okay. I guess, just in that case, what were  
25 they not being -- because I'm not familiar with it,

Page 252

1 paragraph, "In my humble opinion therefore Ryanair is  
2 well within its rights to take this action even though  
3 it is objectionable and puts the consumer in the hot  
4 seat." Do you still hold that opinion today?

5 MR. KASNER: Object to the form.

6 THE WITNESS: I think that if you take the  
7 case of a bad actor, that's if Ryanair was the good  
8 actor here and the person being damaged and the bad  
9 actor was someone like a Bravofly, I'd say yes.

10 BY MS. GIERHART:

11 Q And why is that?

12 A For the reason I think we've already discussed  
13 and for the reasons we have just gone through.

14 MS. GIERHART: All right. I'm going to show  
15 you one more blog post. Let me make sure -- okay.  
16 It's going to be Exhibit 131.

17 (Deposition Exhibit 131 was electronically  
18 received and marked for identification.)

19 BY MS. GIERHART:

20 Q It's from a similar time frame; in fact, I  
21 think only a week later. You can have a minute to read  
22 through it.

23 A I can't see it yet.

24 Q Oh. I'm sorry. Exhibit 131, do you have  
25 that?

Page 251

1 what were they not being clear to the customer about?

2 A I -- again, this is 15-year-old memory, and a  
3 lot of water has flowed under a lot of bridges. I just  
4 remember Bravofly was well known for misrepresenting to  
5 consumers what it did, and it wasn't just an issue of  
6 screen scraping for which they actually did a lot of  
7 screen scraping, not just to Ryanair, but to a bunch of  
8 other players. And it was well known at the time that  
9 Bravofly -- again, this is recalling from memory --  
10 that Bravofly was a bad actor. And I think that when  
11 you look -- put all of those things together, I think  
12 my statement was accurate --

13 Q Okay.

14 A -- for Bravofly.

15 Q And, then, at the end -- okay. Well,  
16 actually, I just want to go to the very end of this,  
17 but to go back up to the start to orient what this  
18 article -- the catalyst for the article, in the second  
19 paragraph, you say this is about the stories that  
20 Ryanair has stated, effectively Monday, they will  
21 automatically cancel any bookings made by screen  
22 scrapers. I think you mean screen scrapers.

23 A I did.

24 Q Okay. So with that context, just to go to the  
25 very bottom, you conclude, the second to the last

Page 253

1 A I don't have it. Oh, sorry. There it is.  
2 Yeah, it came up. Sorry.

3 Q No worries.

4 A Well, this is from August, I believe. The  
5 other one was from June. Was it from June? Oh, no.  
6 This is a week later. You are absolutely right. Okay.  
7 Let me just read this one.

8 Q Sure.

9 A I don't know if I can because it's a pdf. I  
10 was going to say, can we go and look at the article in  
11 Aviation Week? I doubt that we would be able to get to  
12 it --

13 Q Yeah. I'm not --

14 A -- because that would give context.

15 So I think I'm pretty clear here. Again, it's  
16 ten years ago. I was saying the bad Type 2s are  
17 actually the ones causing the problem. I'm trying to  
18 distinguish between bad Type 2s and good Type 2s are  
19 actually the ones causing the problem here.

20 Ryanair is well within the rights to go after  
21 these guys at the time in the jurisdiction that this  
22 happened. If, however, Ryanair went after Type 1  
23 scrapers, then I would be less enthusiastic. By  
24 implication, the good Type 2s would also fall under  
25 that category.

Page 286

1 the same as what a conventional travel agency does.

2 You need to have this minimum information in  
3 order to make a reservation. Yes, that's true. You  
4 can then use that information and look -- when you  
5 combine that with, for example, the behavior of the  
6 customer. So you can get -- you do get personalization  
7 out of this. That's standard business.

8 BY MS. GIERHART:

9 Q Okay. And, then, your next sentence is that,  
10 as such, Ryanair is attempting to establish itself as a  
11 competitor of online travel agencies.

12 Is it your contention that Ryanair is trying  
13 to establish itself as a competitor because it collects  
14 the name, the date of birth, and contact information of  
15 its customers?

16 MR. KASNER: Object to the form. It's one of  
17 the things that Ryanair does. Ryanair wants to provide  
18 the same level of service as an OTA does because, you  
19 know, we've gone through this at great length, looking  
20 at comparing OTAs and direct airlines, and I think  
21 Ryanair is trying to do the same. They want to be  
22 competitive with online travel agencies in gathering  
23 the customers' eyeballs.

24 BY MS. GIERHART:

25 Q You are saying they have the same business

Page 288

1 than what's here in the report?

2 A Sorry. Ask the question again.

3 Q Is there anything else that supports your  
4 contention that Ryanair is attempting to establish  
5 itself as a competitor of OTAs other than what you have  
6 here in the report, which is that they are both  
7 targeting customers?

8 MR. KASNER: Object to the form.

9 THE WITNESS: That is a very broad way that  
10 Ryanair and OTAs are competitive, and they both  
11 consider each other to be competitive. So, in my  
12 report, I've tried to demonstrate how you get to that  
13 point with what is the behavior of an airline, what is  
14 the behavior of intermediaries be they GDS or OTAs or  
15 other forms, and what is the behavior of the customer.

16 So, as a broad brush, I tried to give that --  
17 the answer to that information so that, inside the  
18 report, I tried to give a way that describes how  
19 Ryanair behaves, and that's very similar to how an OTA  
20 behaves, which is what's good.

21 MR. KASNER: I think we are at time. So, if  
22 we could just wrap this up, that would be great.

23 MS. GIERHART: Right. No. Yeah. I think --  
24 I think I'm -- we are at time. So I think I'm done.

25 MR. KASNER: Okay.

Page 287

1 model, which is to collect information about customers  
2 and then target marketing to them; right?

3 MR. KASNER: Object to the form.

4 THE WITNESS: That is one of the things they  
5 do, yes.

6 BY MS. GIERHART:

7 Q Isn't that what essentially every company that  
8 markets to consumers today does?

9 MR. KASNER: Object to the form.

10 THE WITNESS: That's no different. You are  
11 absolutely correct.

12 BY MS. GIERHART:

13 Q So is -- wouldn't that make every company a  
14 competitor of OTAs?

15 MR. KASNER: Object to the form.

16 THE WITNESS: Let's put it this way: Amazon  
17 today -- touch wood -- does not sell travel on  
18 Amazon.com. So I think the context of this is very  
19 important. You can't do a search from Dublin to  
20 Stansted on Amazon.com. You can on an OTA, and you can  
21 on Ryanair.com. So I think the context is very  
22 important.

23 BY MS. GIERHART:

24 Q Okay. Is there anything else supporting your  
25 contention that Ryanair is a competitor with OTAs other

Page 289

1 THE WITNESS: Okay.

2 MS. GIERHART: Did you have any questions, or  
3 are we done?

4 MR. KASNER: We do not have any questions.  
5 We would like to designate the transcript for this as  
6 highly confidential pursuant to the protective order.

7 THE VIDEOGRAPHER: Joanna, do you want copy  
8 orders on the record?

9 THE REPORTER: Yes, please. Does anyone need  
10 a copy of the transcript today?

11 MS. GIERHART: Yes. We'll -- e-tran only. We  
12 don't need paper copies, and if we can, get the rough  
13 and expedited.

14 MR. KASNER: The same for us would be great.

15 THE REPORTER: Thank you.

16 THE VIDEOGRAPHER: All right. This concludes  
17 the deposition on September 26, 2023, at 5:24 p.m.

18 (Deposition session concluded at 5:24 p.m.)

19 -oOo-

Page 290

I certify (or declare) under penalty of perjury under the laws of the State of California that the foregoing is true and correct.

Executed at \_\_\_\_\_ on \_\_\_\_\_.  
(Place) (Date)

\_\_\_\_\_  
(Signature of Deponent)

Page 291

DEPOSITION OFFICER'S CERTIFICATE  
STATE OF CALIFORNIA )  
) ss.

COUNTY OF ORANGE )

I, Joanna B. Brown, hereby certify:  
I am a duly qualified Certified Shorthand Reporter in the State of California, holder of Certificate Number CSR 8570 issued by the Court Reporters Board of California and which is in full force and effect. (Fed. R. Civ. P. 28(a)).

I am authorized to administer oaths or affirmations pursuant to California Code of Civil Procedure, Section 2093(b) and prior to being examined, the witness was first duly sworn by me. (Fed R. Civ. P. 28(a), 30(f)(1)).

I am not a relative or employee or attorney or counsel of any of the parties, nor am I a relative or employee of such attorney or counsel, nor am I financially interested in this action. (Fed R. Civ. P. 28).

I am the deposition officer that stenographically recorded the testimony in the foregoing deposition, and the foregoing transcript is a true record of the testimony given by the witness. (Fed. R. Civ. P. 30(f)(1)).

Page 292

Before completion of the deposition, review of the transcript [ ] was [ ] was not requested. If requested, any changes made by the deponent (and provided to the reporter) during the period allowed, are appended hereto. (Fed. R. Civ. P. 30(e)).

Dated: \_\_\_\_\_

DEPOSITION OF: TIMOTHY JAMES O'NEIL-DUNNE

DATE OF DEPOSITION: Tuesday, September 26, 2023

CASE: *Ryanair DAC v. Booking Holdings Inc., et al*, Case No. 20-1191-WCB**ERRATA SHEET**

The following are the corrections which I have made to my deposition transcript:

<b>Pg.</b>	<b>Ln.</b>	<b>Now Reads</b>	<b>Should Read</b>	<b>Reason</b>
9	12	“(inaudible)”	gist	Transcription
21	21	pr?cis	précis	Transcription
25	14	buying	enabling	Correction
34	11	crew	credentialed	Transcription
58	18	airline	airline,	Transcription
65	19	serve	settle	Correction
81	2	electronic	paper	Correction
102	11	DDS's	GDSs'	Transcription
109	4	obstructed	abstracted	Transcription
112	12	bot	board	Transcription
113	17	Dog	dot	Transcription
141	10	educators	all comers	Transcription
148	9	part	path	Transcription
149	16	office	offers	Transcription
152	7	gain	against	Transcription
153	21	NDCs	NDC	Transcription
156	21	selling	sourcing	Transcription
158	17	agency or a regular brick...	agency or via a regular brick...	Transcription
160	16	Saudi	Saudia	Transcription

DEPOSITION OF: TIMOTHY JAMES O'NEIL-DUNNE

DATE OF DEPOSITION: Tuesday, September 26, 2023

CASE: *Ryanair DAC v. Booking Holdings Inc., et al*, Case No. 20-1191-WCB

Pg.	Ln.	Now Reads	Should Read	Reason
160	16	Saudi	Saudia	Transcription
161	21	indirect	direct	Transcription
167	19	reaching	reach and	Transcription
168	19	Flair,	Flair then,	Transcription
173	19	that	it	Transcription
179	13-14	“— like the”	“— like the – “	Transcription
190	20	airline's website	airlines' websites	Transcription
217	2	Sendai	Cendyne	Transcription
217	3	Sendai	Cendyne	Transcription
219	10	pausing	parsing	Transcription
219	25	pausing	parsing	Transcription
219	15	104	1040	Transcription
220	4	scraping/pausing	scraping/parsing	Transcription
220	21	light	flight	Transcription
223	5	pausing	parsing	Transcription
230	11	DDOS	DDoS	Transcription
232	8	earlier line	Value Alliance	Transcription
232	22	ValueAirlines.com	www.ValueAlliance.com	Transcription
238	19	DDS	DDoS	Transcription
241	6	DDOS	DDoS	Transcription
243	12	DDOS	DDoS	Transcription
244	23	DDOS	DDoS	Transcription

DEPOSITION OF: TIMOTHY JAMES O'NEIL-DUNNE

DATE OF DEPOSITION: Tuesday, September 26, 2023

CASE: *Ryanair DAC v. Booking Holdings Inc., et al*, Case No. 20-1191-WCB

Pg.	Ln.	Now Reads	Should Read	Reason
260	19	customer's water	customer is like water	Transcription
261	21	FAA	2FA	Transcription
267	21	pr?cis	précis	Transcription
277	24	DDOS	DDoS	Transcription
286	16	form. It's one	form." [Insert new line] "THE WITNESS: It's...	Transcription. The transcript improperly combines Mr. Kasner's objection with The Witness's answer
293	Index	DDOS	DDoS	Transcription
307	index	DDS's	GDSs'	Transcription
307	Index	DDS	DDoS	Transcription
307	Index	DDOS	DDoS	Transcription
331	Index	pausing	parsing	Transcription
334	Index	pr?cis	précis	Index should be corrected to reflect "précis"
340	Index	scraping/pausing	scraping/parsing	Transcription

I, the undersigned, declare under penalty of perjury, that I have read the above-referenced deposition transcript and have made any corrections, additions or deletions reflecting my true and correct testimony.

EXECUTED this 15 day of November 2023, at 10:15 P<.

Timothy James O'Neil-Dunne  
Timothy James O'Neil-Dunne

# EXHIBIT 13

**PUBLIC VERSION -  
CONFIDENTIAL MATERIAL OMITTED IN FULL**



# EXHIBIT 14

**PUBLIC VERSION -  
CONFIDENTIAL MATERIAL OMITTED**

# Web Scraping for Hospitality Research: Overview, Opportunities, and Implications

Cornell Hospitality Quarterly  
2021, Vol. 62(1) 89–104  
© The Author(s) 2020  
Article reuse guidelines:  
sagepub.com/journals-permissions  
DOI: 10.1177/1938965520973587  
journals.sagepub.com/home/cqx



Saram Han<sup>1,2</sup>  and Christopher K. Anderson<sup>1</sup> 

## Abstract

As consumers increasingly research and purchase hospitality and travel services online, new research opportunities have become available to hospitality academics. There is a growing interest in understanding the online travel marketplace among hospitality researchers. Although many researchers have attempted to better understand the online travel market through the use of analytical models, experiments, or survey collection, these studies often fail to capture the full complexity of the market. Academics often rely upon survey data or experiments owing to their ease of collection or potentially to the difficulty in assembling online data. In this study, we hope to equip hospitality researchers with the tools and methods to augment their traditional data sources with the readily available data that consumers use to make their travel choices. In this article, we provide a guideline (and Python code) for how to best collect/scrape publicly available online hotel data. We focus on the collection of online data across numerous platforms, including online travel agents, review sites, and hotel brand sites. We outline some exciting possibilities regarding how these data sources might be utilized, as well as discuss some of the caveats that have to be considered when analyzing online data.

## Keywords

web scraping; online review; data collection; Python

## Introduction

The internet presents many interesting opportunities for understanding consumer choice. Any information displayed on any travel website represents a potential data source for hospitality researchers. For instance, a hospitality researcher might collect a list of reviews from TripAdvisor and perform text mining, or retrieve satisfaction ratings to perform some statistical analysis, or determine the relationship between price and page ranking at an online travel agent (OTA). Despite the appeal of these opportunities, however, the fact remains that manually collecting a large amount of online data is inefficient and practically impossible. A typical hospitality research study may require data from multiple hotels across different markets over a prolonged period of time; collecting this data manually can be incredibly time consuming and tedious. Therefore, hospitality researchers need to automate the process of gathering and storing the information presented on travel websites. This is where web scraping comes into play. Web scraping is the process of creating a computer program to download, parse, and organize data from the web in an automated manner (vanden Broucke & Baesens, 2018).

As Table 1 illustrates, there are several alternative approaches that hospitality researchers may consider when trying to automate the collection of online data. One

promising option is using an Application Programming Interface (API) provided by the firm hosting the data. However, APIs are often difficult to access, accessible only for a limited duration, or require upfront fees even when researchers are able to access an API. Furthermore, APIs generally will not expose all the required variables.

Recently, scholars have utilized alternative approaches to collect consumer behavior data that mimics online behaviors. Previous studies directly captured data in a simulation lab or survey setting where research participants acted as if they were actually completing online activities (Choi et al., 2017; Min et al., 2015; Shin et al., 2019). However, a common criticism of these studies is that there is a mismatch between real customer behavior and the behavior of the subjects in these experiments, which are often conducted by undergraduate students or participants in Amazon's Mechanical Turk (Schahn & Holzer, 1990). Moreover, although it is possible to conduct large-scale

<sup>1</sup>Cornell University, Ithaca, NY, USA

<sup>2</sup>Seoul National University of Science and Technology, Republic of Korea

### Corresponding Author:

Saram Han, Seoul National University of Science and Technology,  
232 Gongneung-ro, Nowon-gu, Seoul 01811, Republic of Korea  
Email: saramhan@seoultech.ac.kr

**Table 1.**  
**Comparison of Common Data Collection Methods in Online Market Research.**

	Scraped Data	Commercial Web Scraping Service	API	Survey
Cost	Low	Medium	Low/Medium	High
Sample frame	Website users	Website users	Website users	Flexible
Customizability of variables	Medium	Low	Low	High
Ease of frequent collection	Easy	Moderate	Easy	Hard
Data type	Behavioral	Behavioral	Behavioral	Attitudinal
Limitations	Time and programming skills	Data may not be suitable to the researcher's need in terms of variables or content	Limited availability	Time and programming skills

Note. API = Application Programming Interface.

survey research, a number of barriers make that process difficult. Meanwhile, the scale of web scraped data usually is large in nature. Compared with conducting a survey, collecting publicly available online data is inexpensive—even considering the required fee for the commercial service. The monetary cost and effort required to conduct a survey increases linearly with each additional sample, whereas web scraping requires only a one-time expense. Hotel prices change daily, and online reviews are posted even more frequently. In such a dynamic market, a single sampling may not be sufficient. Conducting multiple surveys sequentially is expensive and requires a lot of effort, whereas collecting publicly available online data enables researchers to extract as much data as they need in real time.

Another data collection method that researchers frequently chose is to outsource the scraping task to third-party commercial firms (Wu et al., 2015). Although this is perhaps the simplest option, there are potential pitfalls that may arise due to the lack of discretion in the data collection process. For instance, researchers often realize too late that they forgot to ask the firm to scrape for certain variables. As a result, they are forced to either spend additional time, effort, and money or complete their research without having access to all the desired variables or behaviors.

However, the existing literature points out several drawbacks to web scraped data. One common critique of online data is that study results derived from this data are rarely generalizable to the behavior of the entire target population; although these data are useful to researchers who are interested only in the behavior of online customers, if the researcher's target population includes both online and offline customers, any results derived from online data will be highly vulnerable to selection bias. This is because web scraped data may overrepresent the unique behaviors of online customers. Meanwhile, researchers who opt to utilize surveys as their research method have the flexibility to design a sample frame capable of representing both online and offline customers. Another drawback of publicly

available data is that researchers have little discretion in measuring aspects from the sample. In contrast, researchers can include anything they want in a survey or questionnaire, and thus have no restriction in the variables that they can collect. Researchers should carefully decide whether web scraping is suitable to their research and fully consider the potential challenges web scraping poses. An ideal approach might be to collect data through both experiments and web scraping to confirm both the internal and external validity (Viglia & Dolnicar, 2020). For example, some researchers use web scraped data to develop their hypotheses and use traditional data collection methods to confirm these hypotheses (Kupor & Tormala, 2018). Although web scraping has become popular recently and presents an important opportunity to better understand the online marketplace, traditional data collection methods remain the most commonly used strategies in hospitality research.

Table 2 illustrates the potential benefits of web scraped data for hospitality research by summarizing papers by data type with a focus on online reviews in travel marketplaces. To gather this data, we surveyed six highly reputable hospitality journals using the advanced search feature that is available on any journal website. Our analysis included all papers written prior to July 2020. We classified papers as having utilized traditional data collections methods (such as surveys or interviews) if they included the term “online review” and at least one of the terms {“interview,” “survey,” “amazon mechanical turk,” “questionnaire”}, but none of the terms {“scrape,” “scraping,” “crawl,” “actual review,” “python”}. In contrast, any papers including the term “online review” and at least one of the following terms {“scrape,” “scraping,” “crawl,” “actual review,” “python”} were classified as having directly used web data for online review research. The search queries we used are as follows:

- Interview or survey: “online review” AND (“interview” OR “survey” OR “amazon mechanical turk” OR “questionnaire”) NOT (“scrape” OR “scraping” OR “crawl” OR “actual review” OR “python”).

**Table 2.**  
**Data Collection Methods Used for Online Review Research in Hospitality Journals.**

	Interview OR Survey	Web Data
<i>Cornell Hospitality Quarterly</i>	16	3
<i>Journal of Travel Research</i>	21	4
<i>Journal of Hospitality and Tourism Research</i>	7	3
<i>International Journal of Hospitality Management</i>	138	21
<i>Annals of Tourism Research</i>	31	5
<i>Tourism Management</i>	82	25

- Web data: “online review” AND (“scrape” OR “scraping” OR “crawl” OR “actual review” OR “python”).

As the table indicates, only a handful of papers use data from the internet to tackle research questions that are related to online marketplaces.

The increasing availability of open source web scraping tools has made it a lot easier for the researchers to build their own customized web scrapers. This article aims to dismantle some of the barriers that hospitality researchers encounter when attempting to utilize web scraping in their online studies. However, this article is not a complete guide to web scraping, but rather an introduction to some of the key tools and requirements for scraping key hospitality data sources. Our article provides hospitality researchers with tools and techniques for collecting data from typical, interactive hotel websites. Due to the introductory nature of this article, we focus on applied concepts and functional illustrations. In this article, we assume that readers have fundamental knowledge of programming languages, such as defining variables, loops, and functions, but not necessarily Python. Proficient programmers and researchers who are familiar with web scraping may already be capable of writing their own code from scratch, and thus may find this article to be too applied or specific. The methods we suggest may not represent the most efficient methods available, as this article only considers websites that have dynamic contents—that is, webpages designed to change in response to human interactions (Massimino, 2016): for example, say a researcher wants to scrape not only all the customer reviews on a particular hotel webpage on TripAdvisor but also the profile information from each individual reviewer. TripAdvisor’s dynamic website is designed to present detailed profile information only when a user positions their cursor over the reviewer’s profile picture. This is where the dynamic-website specialized code that we cover in this article becomes useful. In contrast, if the website contains only static contents—meaning that the website does not change until the user (i.e., client) moves on to another webpage—our dynamic-website specialized code may take longer to function than static-website specialized code. Therefore, for

efficient coding, some researchers may want to refer to Massimino (2016) if the website contains only static contents. While the dynamic-website specialized codes can be used to analyze both static and dynamic contents, static-website specialized codes cannot analyze dynamic contents. As the goal of this article is to provide hospitality researchers with a generally useful tool that can be easily applied to a variety of websites, we only consider the dynamic-website specialized code in this article. Popular Python libraries specialized for static website scraping are Requests, BeautifulSoup, lxml, and Scrapy. Those who are interested in detailed instructions from utilizing these libraries and other web scraping methods that can be applied to broader websites should read vanden Broucke and Baesens (2018). We hope our article will make online data collection easier for hospitality researchers and spark greater interest in web scraping techniques. Therefore, our purposes include the following:

- Providing tools applicable to major hotel platforms (i.e., TripAdvisor, Expedia, Marriott.com, and Airbnb),
- Introducing the process in the simplest terms possible, and
- Discussing both the benefits and limitations of analyzing web scraped data.

This article is a step-by-step guide to web scraping, with a focus on websites that are particularly useful to hospitality researchers. Other papers tend to either focus too narrowly on specific websites or discuss web scraping at an abstract level. Considering that different travel websites provide different insights, hospitality researchers rarely utilize only one specific website in their research. Narrowly focused articles require hospitality researchers to learn different web scraping methods for each individual website included in their research. However, overly abstract articles often fail to provide all-in-one, generalized information that novices can apply to their own research. To narrow the gap between the needs of hospitality researchers and the currently available resources, our article focuses on two aspects that have not previously been thoroughly examined in a single paper.

First, rather than introducing different tools for each individual website, we introduce more general tools that are applicable to all major travel websites. Second, after reading our article, researchers will have all the knowledge they need to start web scraping, as our article provides the entire scripts used for scraping every major travel website.

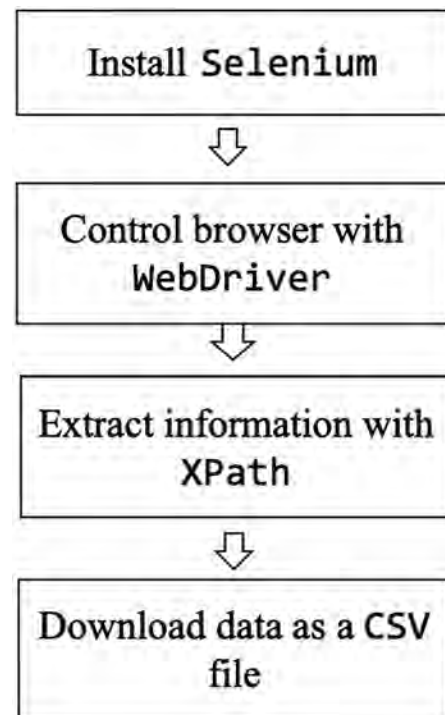
The remainder of this article is constructed as follows. First, we introduce the environments required for running a Python-based scraper. Then, we explain the sample web scraping code for collecting online hotel reviews and prices. In addition to providing samples of key code required for scraping through the paper, in the supplemental appendix we provide complete scripts for scraping reviews and prices from major OTAs, review sites and hotel brand sites. As most web scraped data are secondary data, with the data generation process outside the researcher's control, biases have to be handled with caution. Therefore, we also illustrate some possible biases that researchers should be aware of before analyzing scraped data. Finally, we discuss the academic implications of web scraping and the ethical issues that hospitality researchers should consider when collecting online data for their own research.

## Web Scraping Fundamentals

In the following section, we outline the fundamentals of web scraping illustrated through the use of Python—a general-purpose programming language. Our article is not an introduction to Python as a whole, but rather focuses on the aspects of Python that are key to web scraping. Python is a very approachable and intuitive programming language and is often the language of choice for many introductory programming courses. There are many online resources that provide an overview of Python. For those who are interested in learning the basics of Python language, we recommend reading Downey (2014). However, our focus in this section is on the unique methods and skills necessary for scraping hospitality data from a variety of sources using Python. As there is a large scraping community that uses Python3, throughout the article we utilize Python3 specifically (<https://www.python.org/downloads/>).

### Required Environment and Python Codes

As our goal is to simplify the process of gathering online information and transforming it into a meaningful data set, we focus only on the essential elements of web scraping in this article. Over the course of this article, we use three different tools to make a scraper that facilitates the collection of data from major hotel websites: Selenium, WebDriver, and XPath. Figure 1 presents the flowchart of the steps of how we utilize the tools specialized for web scraping the dynamic contents.



**Figure 1.**  
Flowchart of Web Scraping for Extracting Dynamic Contents.

**Selenium.** In this section, we introduce the Python library Selenium. Although there are many libraries in Python that facilitate web scraping, Selenium is the most useful when dealing with interactive, JavaScript-heavy pages like those on travel sites such as TripAdvisor, Airbnb, and Expedia. We start by illustrating scraping approaches for extracting review data from TripAdvisor (Listing 1). We then extend this approach to extract review data from other platforms and later extend this approach further to collect other types of data, for example, prices. Let's begin by scraping reviews and basic reviewer profile information from the following TripAdvisor link: [https://www.tripadvisor.com/Hotel\\_Review-g60763-d93344-Reviews-The\\_Watson\\_Hotel-New\\_York\\_City-New\\_York.html](https://www.tripadvisor.com/Hotel_Review-g60763-d93344-Reviews-The_Watson_Hotel-New_York_City-New_York.html).

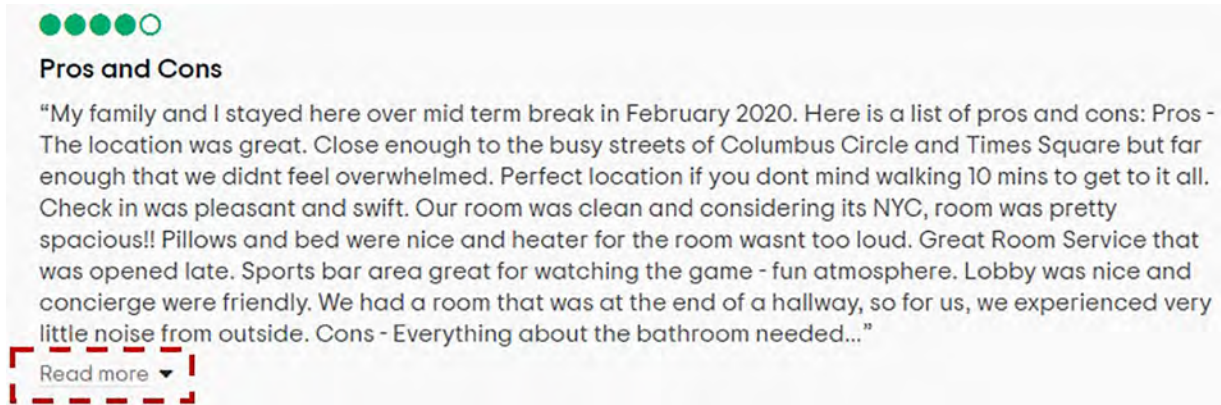
### Listing 1. Define Our Target Website for Scraping.

```

1. target_page = "https://www.tripadvisor.com/Hotel_Review-g60763-d93344-Reviews-The_Watson_Hotel-New_York_City-New_York.html"
  
```

First, let's name this target page. With other Python libraries (BeautifulSoup, request, Scrapy, etc.), you can





**Figure 2.**  
Example of an Interactive Platform.

scrape everything that is displayed on the screen. While this is good news, it also means that you cannot scrape information that is hidden until you click on it. For instance, contents that are hidden behind the button “Read more,” as illustrated in Figure 2, cannot be scraped with other popular Python libraries such as BeautifulSoup, request, or Scrapy. Unlike websites that are only written in Hypertext Markup Language (HTML) or CSS, many interactive travel sites use JavaScript, which allow for interactive functionality on the web page. Contrary to many other programming languages, the core functionality of JavaScript lies in making webpages more interactive and dynamic. Therefore, for sites that make heavy use of JavaScript, we need to write our code in a way that emulates human browsing behavior. That is where the Python library Selenium comes into play. You can easily install it with the following code in the command line: `pip install -U selenium` from the terminal window on your computer. Selenium requires a third-party software called a WebDriver which we discuss in the next paragraph.

**WebDriver.** WebDrivers exist for most modern internet browsers, including Chrome, Firefox, Safari, and Internet Explorer. When using these browsers, a browser window will open up on your screen and perform the actions specified in your code. We can easily download a WebDriver from <https://sites.google.com/a/chromium.org/chrome-driver/downloads>. As WebDrivers exist for most popular browsers, you can choose to download whichever WebDriver best works for you. In the following section, we assume that readers will download the WebDriver for Google’s Chrome browser—chromedriver. Make sure you download the file that matches both your operating system (i.e., Windows, Mac, or Linux) and the version of your current browser (i.e., Chrome, Internet Explorer, Edge, or Firefox). The ZIP file you download will contain an executable called `chromedriver.exe` on Windows, or simply

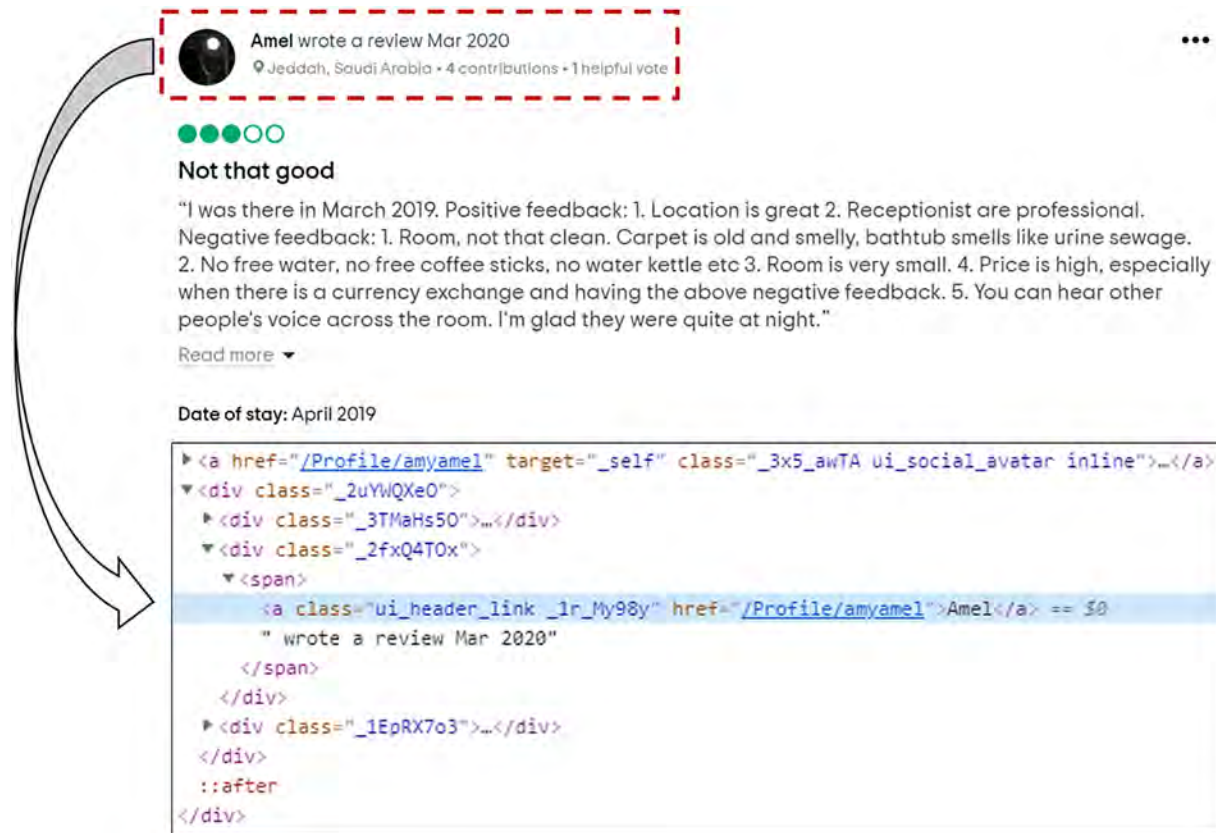
`chromedriver` otherwise. While locating the WebDriver in the same directory as your Python scripts is the easiest way to call the driver, it is also possible to explicitly pass the location of where you have the WebDriver as we do with the variable `driver_loc` in the code Listing 2. Once the WebDriver opens the browser, you can interact with the browser until you close it (i.e., `driver.close()`).

**Listing 2.**  
Python Code for Opening and Closing Your WebDriver.

```
1. driver_loc = 'C:/Users/Desktop/chrome
   driver.exe'
2. driver = webdriver.Chrome(driver_loc)
3. driver.get(target_page)
4. #Write your scraper here.
5. driver.close()
```

**XPath.** Once a page is loaded with WebDriver, you will want to extract the information displayed on the screen and download it on your computer. To extract this information, you first need to view the source code of the website elements you want to scrape. For Windows users, the easiest way to check this information is to right-click on the part within the web page you want to scrape and select “Inspect.” You will then see a tree-based view of the HTML codes. Every web page is made up of a bunch of these HTML tags denoting each type of content on the page. By clicking on the arrows, you can see the nested structure of the code. To scrape hotel review websites, we only need to understand the basic structure of HTML language, which consists of tag, attribute, and the value of the attribute. Figure 3 illustrates part of the source code underlying the page displayed after selecting “Inspect.”

As an illustration, let us assume that the element of the HTML code written to display the reviewer’s name (i.e.,



**Figure 3.**  
Displayed Tree-Based HTML Code View after Clicking on Inspect.

Amel) is shown in Listing 3, where the `<a> . . . </a>` is just one of the many existing tags (e.g., `<p>`, `<div>`, `<ul>`) in HTML which encloses the hyperlink and styles the web designer decided to use when displaying the reviewer's name. The value of the attribute "MemberBlock" is the name of a style that is assigned to the attribute class. Likewise, href is another attribute of this element that redirects users to the reviewer's profile page when they click on the reviewer's username. XPath is a language used to find the location of any element on a web page using this structure of the HTML. Selenium can use the XPath language to select elements. Although there are many Selenium methods, we use only `find_element_by_xpath` and `find_elements_by_xpath` throughout this article for ease of application.

**Listing 3.**  
HTML Element for the Reviewer's Name.

```
1. <a class="MemberBlock" href="/Profile/amyamel">Amel</a>
```

To extract the information from the HTML code Listing 3, the XPath has to be written as, `//a[@`

`class='MemberBlock']`, which basically is `//Tag[@Attribute='Value']` that can be interpreted as "grab the current node (//) with the tag name where the class attribute corresponds to MemberBlock." However, oftentimes the value of the attribute is very long and contains multiple special characters. This can be difficult because if even one of the characters is missing, a serious error may occur. `//a[contains(@class, 'Member')]` presents an easier way to call the same element, by utilizing the same process but not requiring users to include the entire value of the attribute. This XPath prevents the researcher from provoking an error simply by not indicating the exact value of the attribute. The function `contains()` enables Selenium to easily find the desired element, even if the researcher only provides a fragment of the attribute value. Accordingly, the `find_element_by_xpath` method in the code Listing 4 indicates the element of the HTML code Listing 3.

**Listing 4.**  
Indicating the HTML Element.

```
1. driver.find_element_by_xpath("//a[contains(@class, 'Member')])")
```



To extract the reviewer's name, we can add `.text` as shown in the code Listing 5.

**Listing 5.**  
**Indicating the Text Content.**

```
1. driver.find_elements_by_xpath(".//a
   [contains(@class, 'Member')]").text
```

Note that the value of the same attribute can always change (e.g., `MemberBlock`) whenever the web designer defines a different name for the value. Therefore, the same web scraper may not work if the web designer updates the HTML code. If an error occurs after running the codes that we provided in the supplemental appendix, we recommend our readers to review the HTML code of the target website and check whether the value of the attribute has changed.

If the researcher is interested in extracting the user's profile link, the following `.get_attribute("href")` method is useful for finding the value of the `href` attribute (Listing 6).

**Listing 6.**  
**Indicating the Value of the "href" Attribute.**

```
1. driver.find_element_by_xpath(".//a
   [contains(@class, 'Member')]").text.
   get_attribute("href")
```

### Applications for Major Hotel Platforms

In this section, we apply these scripts to scrape four major hotel review websites and discuss how to best deal with platform-specific features when building a scraper. We also discuss how hospitality researchers can take advantage of the unique opportunities presented by each individual platform.

**Reviews on TripAdvisor.** TripAdvisor is the largest player in the travel review platforms arena (Wang & Chaudhry, 2018). Therefore, this platform has been used in many marketing and hospitality studies (Chevalier et al., 2018; Wang & Chaudhry, 2018). TripAdvisor has several unique features that make it attractive to researchers. For instance, each individual's previous platform activities are displayed on their profile. This feature allows researchers to take reviewer level heterogeneity into account as they analyze TripAdvisor data (Gao et al., 2017). Therefore, we recommend scraping each reviewer's profile URL to collect individual-level data that may be useful in the future.

Another interesting feature that recent studies have started to focus on is the existence of *retailer-prompted reviews* (Askalidis et al., 2017; Han & Anderson, 2020;

Mayzlin et al., 2014). Unlike regular self-motivated reviews, retailer-prompted reviews are those that are posted in response to hotels' email invitations to post reviews online. This feature allows researchers to investigate how rating behavior differs depending on the nature of the review posting process.

Although most scraping tasks can be completed using the aforementioned functions and XPath, there is one problem that requires additional attention. As mentioned in Figure 2, many TripAdvisor reviews are long, meaning that users must click on the "Read more" button to see the full review. Researchers must be able to automatically click this button to scrape entire reviews. For this purpose, Selenium offers a selection of "actions" that can be performed by the browser, such as clicking elements. Fortunately, in TripAdvisor, once you expand one review, the rest of the reviews within the same page expand as well. If for some reason all of the desired reviews do not automatically expand with a single click, a for-loop can be added to expand each review individually (see Listing 8). Once we identify the element where the "Read more" button locates, we can use the `execute_script` method to expand the review. We recommend using try-except to avoid the scraper stopping in situations where there is no expandable review on the page.

**Listing 7.**  
**Python Script to Expand the Review Element.**

```
1. more = driver.find_elements_by_xpath
   ("//div[contains(@data-test-target,
   'expand')]")
2. try:
3.     driver.execute_script("arguments[0].
   click();", more[0])
4. except:
5.     pass
```

**Listing 8.**  
**Python Script to Expand Every Review Element Within a for-loop.**

```
1. for container in containers:
2.     more = container.find_elements_by_
   xpath(".//div[contains(@class,
   'expand')]")
3.     try:
4.         driver.execute_script("arguments
   [0].click();", more[0])
5.     except:
6.         pass
```



**Figure 4.**  
Pagination.

More reviews

**Figure 5.**  
Infinite Scrolling.

TripAdvisor presents hotel reviews over multiple separate pages. Once you have successfully scraped the information from the first page of reviews, you may want to move on to the following page to scrape older reviews. This is called pagination (Figure 4; Zhang et al., 2020). We can use the `execute_script` method again to send a JavaScript command to the browser. This method clicks on the “Next page” button, enabling researchers to scrape the same information on the following pages (Listing 9).

**Listing 9.**  
Clicking the Next Page Button.

```
1. element = driver.find_element_by_xpath('//a
   [contains(@class, "nav next")])
2. driver.execute_script("arguments[0].
   click();", element)
```

**Reviews on Expedia.** One big advantage of Expedia is that all reviews are written by valid customers. That is, as Expedia only allows customers who purchased through their website to write reviews, there is a smaller chance that users will see fake reviews written by non-verified customers on Expedia than on other websites such as TripAdvisor, where anyone can write reviews. Moreover, as Expedia sends out email requests to everyone who makes purchases through their website, the entire review collection process is similar to a survey. Therefore, as long as we know who did not respond to the email request sent by Expedia, we can generalize the study results to the population of those who purchased through Expedia. This may not be possible when using data from TripAdvisor, as the TripAdvisor customers’ sample frame is unknown.

On Expedia, users may not need to click the “Next page” button as they do on TripAdvisor. Instead, they must scroll down and load older reviews by clicking on the “More reviews” button shown in Figure 5. This is called infinite scrolling. If there is a significant quantity of reviews written

about a given hotel, users must scroll further down and repeatedly click on the “More reviews” button until all the reviews have loaded. The scraper must mimic this process. Listing 10 opens the number of reviews that you assign with `revnum` (e.g., we set here as 300). This while-loop continues scrolling down until the accumulated number of loaded reviews exceeds `revnum`.

**Listing 10.**  
Scrolling Down.

```
1. revnum = 300
2. loadednum = 0
3. while loadednum < revnum:
4.     more = driver.find_element_by_xpath(
       ("//button[contains(@class, 'more-
       reviews-button')]")
5.     driver.execute_script("arguments[0].
       click();", more)
6.     time.sleep(1)
7.     containers = driver.find_elements_
       by_xpath("//div[contains(@class,
       'uitk-card-separator-bottom')]")
8.     loadednum = len(containers)
```

**Reviews on Airbnb.** Airbnb is a peer-to-peer marketplace that emerged as a typical example of what is called the sharing economy. Many studies present evidence that Airbnb threatens the traditional accommodation market system (Zervas et al., 2018). This website is most interesting to researchers hoping to better understand customers’ behavior in the sharing economy. Due to the peer-to-peer market nature, previous studies argue that there are unique rating behaviors in Airbnb (Ert & Fleischer, 2019). For instance, in addition to guests rating service providers, guests are also evaluated by the host, and these ratings are made publicly available as well. Owing to this dual review process, suppliers tend to receive overwhelmingly high ratings that are not observed under other hotel review systems (Zervas et al., 2018). Despite these unique features of the sharing economy, little is understood about this system. We hope that, by making Airbnb data more accessible, we can encourage other researchers to further explore the sharing economy.

The scraping code for Airbnb is a combination of the codes used for TripAdvisor and Expedia, as the website randomly changes their review display. That is, when the WebDriver visits the supplier's page on Airbnb, the reviews are randomly listed either as pagination or infinite scrolling. Therefore, once the WebDriver opens Airbnb, we need to first determine which scraper should be applied based on the HTML structure. Although there are many different ways to accomplish this, here we build a `check_exists` function that checks whether Airbnb reviews are displayed in pagination or infinite scrolling.

```
1. def check_exists(self, xpath):
2.     try:
3.         self.find_element_by_xpath(
4.             xpath)
5.     except NoSuchElementException:
6.         return False
```

After building the function that checks which platform design we have to deal with, we can simply run different scrapers by writing the following if statement.

```
1. # example of xpath that uniquely
   indicates Infinite scrolling
2. xpath = '//*[@class="_16i7snfh"]'
3. if check_exists(driver, xpath)==True:
4.     # Codes for infinite scrolling
5. else:
6.     # Codes for pagination
```

**Reviews on hotel brand sites.** Major hotel brands have their own websites where customers can write reviews about their experience, just as they can on TripAdvisor or Expedia. Brands offer up their own reviews in an effort to reduce the need for prospective consumers to visit other travel sites such as TripAdvisor or Expedia. Hotel brands are motivated to do this because reservations made on sites such as TripAdvisor and Expedia are more costly to hotels and because customers who visit these sites may end up booking their stay with another company. It is possible that customers who read and write reviews on hotel brand websites differ from those who use other platforms where alternative hotel options are listed. For instance, these customers may be more loyal to the brand or hoping to have their opinions heard by the hotel manager rather than by other potential customers. Platforms like this, where the user groups differ from other platforms users, generate lots of interesting research opportunities.

Although the scraper should be written differently depending on whether the targeted hotel website uses pagination or ultimate scrolling, the overall process of writing the scraper remains exactly the same regardless of the target website. We include an illustration of Marriott.com in Appendix A (along with complete scripts for other travel sites).

### Scraping (Prices) at the Market-Level

While scraping online reviews enable researchers to understand customer Word-of-Mouth (WOM) behaviors, scraping prices gives us insight into the market. This requires no major additional functions or methods beyond the scripts that we introduced previously for scraping reviews. The only difference is that we loop over different hotels within a specific market instead of different reviews. Remember that our scraping code differs by whether the page uses pagination or infinite scrolling. While hotel prices listed on Airbnb and TripAdvisor are displayed using pagination, hotels on Expedia are listed using infinite scrolling. Therefore, utilizing the same approach detailed in code Listing 9, we can build our scraper for Airbnb and TripAdvisor to collect price data and click on the "Next" button to move on to the following pages (see code Listing 11).

#### Listing 11.

**Click the Next Button to Collect all Hotel Prices.**

```
1. next_exists = check_exists(driver, '//a
   [contains(@aria-label, "Next")]')
2. if next_exists:
3.     element=driver.find_element_by_
4.         xpath('//a[contains(@aria-label,
5.             "Next")]')
6.     driver.execute_script("arguments[0].
7.         click();", element)
```

In contrast, for Expedia we design the code to scroll down until there are no additional listings before documenting the hotel prices (see code Listing 11). This approach is equivalent to how we scraped the reviews from Expedia in code Listing 12.

#### Listing 12.

**Scrolling Down to Collect all Hotel Prices.**

```
1. while check_exists(driver, "//button
   [contains(@data-stid, 'show-more-
   -results')]"):
2.     more = driver.find_element_by_xpath
3.         ("//button[contains(@datastid,
4.             'show-more-results')]")
```

**Table 3.**  
The Number and the Average Length of Reviews Across Three Platforms During 2018–2020.

	TripAdvisor	Expedia	Hotel Brand
Number of reviews	713	2,089	1,760
Average review length	731.12	119.56	141.63

```

3.     driver.execute_script("arguments[0].
      click();", more)
4.     time.sleep(1)
5.     containers = driver.find_elements_
      by_xpath("./div[contains(@class,
      'link-container')]")

```

### Platform Differences Causing Biases

Although collecting and analyzing online review data broadens our understanding, it is important to mention a few relevant caveats. Unfortunately, many studies conducting research using data from online travel sites ignore the biases that arise when using secondary data. Multiple selection biases exist in scraped data, which can impact a researcher's ability to draw insights about the target population (i.e., all customers who stayed at the hotel). In this section, we present some of the noteworthy differences between hotel review platforms that hospitality researchers should be aware of before collecting and using web scraped data.

### Review Differences Across Platforms

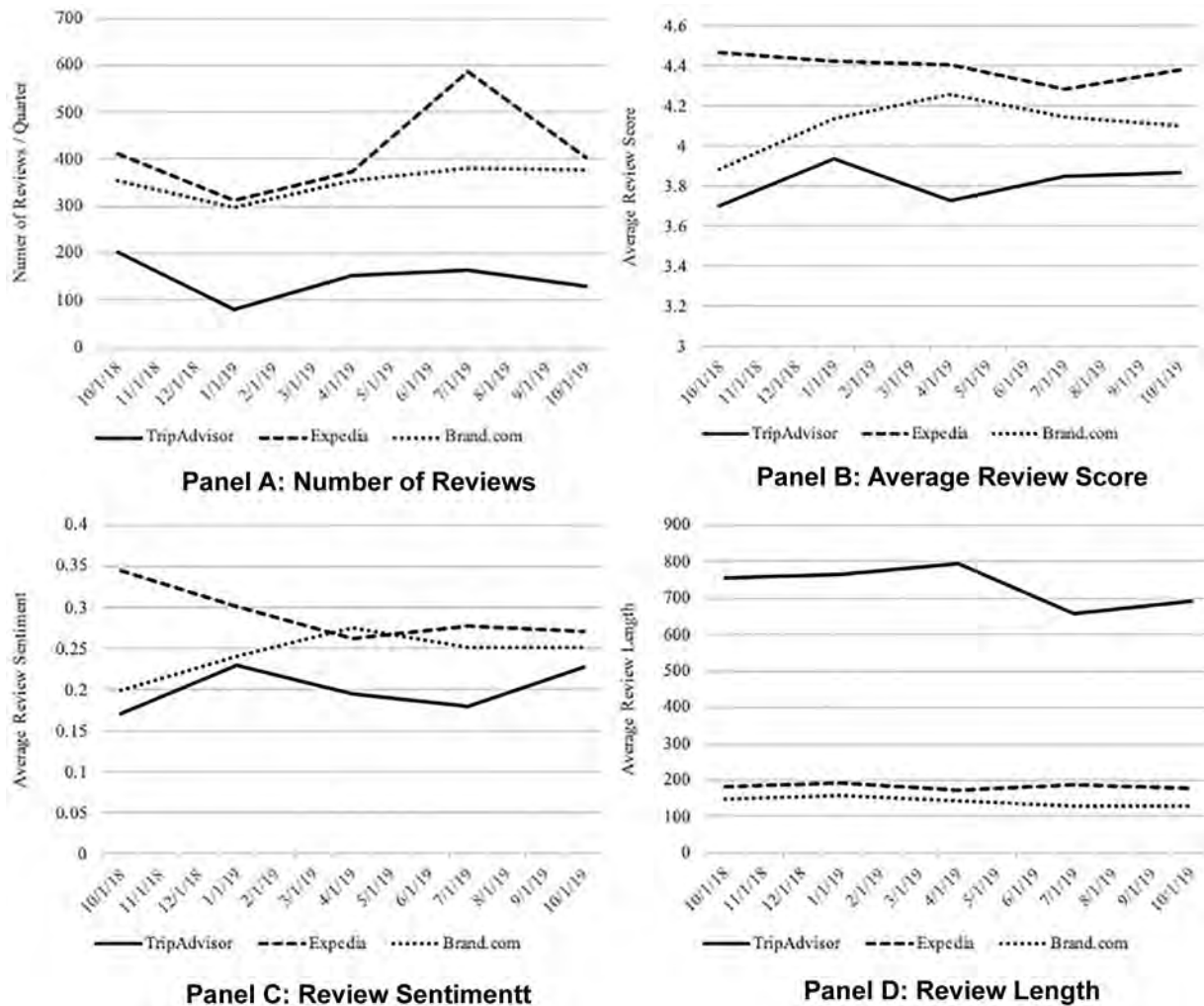
Owing to differences in how reviews are collected and written, collected reviews may vary in terms of number, valence, and content depending on what hotel review platform they were posted on (Litvin & Sobel, 2019). This is true even for reviews that were written about the same hotel and during the same time period. This is true for a couple of reasons. First, every hotel review platform has its own objectives and its own review collection process. To encourage review submission, Expedia sends customers a post-stay email with a link to submit a hotel review. TripAdvisor, however, relies fully on customers' self-motivation to post reviews: unless individual hotels choose to partner with TripAdvisor to invite their verified customers to post reviews, either by sending an email or through their online reputation management firms (e.g., ReviewPro, Revinate, and Medallia). These differences in the review collection process among these platforms yield systematic differences in terms of the review scores and their content. Second, customers may perceive different platforms as having different audiences. Therefore, different types of customers may prefer different platforms depending on who they hope to reach with their review. While customers may perceive distributors' websites (i.e., a hotel's own website or Expedia) as a suitable

channel for contacting hotel managers, TripAdvisor might be perceived as an appropriate method for reaching other potential customers. As a result, the dominant valence and topics may differ across platforms.

Several noteworthy patterns are frequently observed across travel platforms. First, a major difference between review platforms is the number of reviews per hotel. To illustrate how a given hotel's quality might be evaluated differently by each of the three platforms, we scraped the reviews of the New York Marriott Marquis, which has a large number of reviews on TripAdvisor, Expedia, and Marriott's own website. For the purposes of this demonstration, we utilized reviews written between October 2018 and December 2019. Table 3 demonstrates that significantly more reviews were posted on Expedia and the hotel's own website than on TripAdvisor. This pattern is particularly interesting because anyone can post a review on TripAdvisor, whereas only verified guests can post reviews on OTAs and hotel brand sites. A major driver in creating this disparity in the number of reviews across the platforms stems from how the reviews are collected. Expedia and hotels encourage customers to post reviews online by sending post-stay feedback invitation emails. This pattern holds across different time periods, as shown in Figure 6, where we plot review characteristics (number, score, sentiment, and length) on a quarterly basis for 2018 and 2019. The average review length is significantly longer on TripAdvisor than on the other platforms. Self-motivated (i.e., non-email prompted) reviewers are more likely to put greater effort into writing reviews, given that they have already undertaken the effort of opening TripAdvisor and posting a review. Therefore, review length (calculated in number of words) is greatest on TripAdvisor, where most reviews are self-motivated. After eliminating stop words, which refer to the words that are filtered out due to their extremely high or low appearances, we counted the words that appeared on each platform. As is shown in Table 4, the most frequently used words are slightly different depending on the platform. Unlike Expedia and TripAdvisor, where the most frequent word is "room," the hotel's own site reviewers are more focused on the hotel location.

Finally, the distribution of ratings between the two platforms has different patterns, as is shown in Figure 7. The proportion of relatively negative ratings is higher on TripAdvisor than on Expedia or the brand site. Given the cost or effort of posting a review on TripAdvisor, customers may be more likely to post if their expectations have not been met (or have been exceeded), and as a result these





**Figure 6.**  
Differences in Reviews (Number, Score, Sentiment, and Length) by Platform.

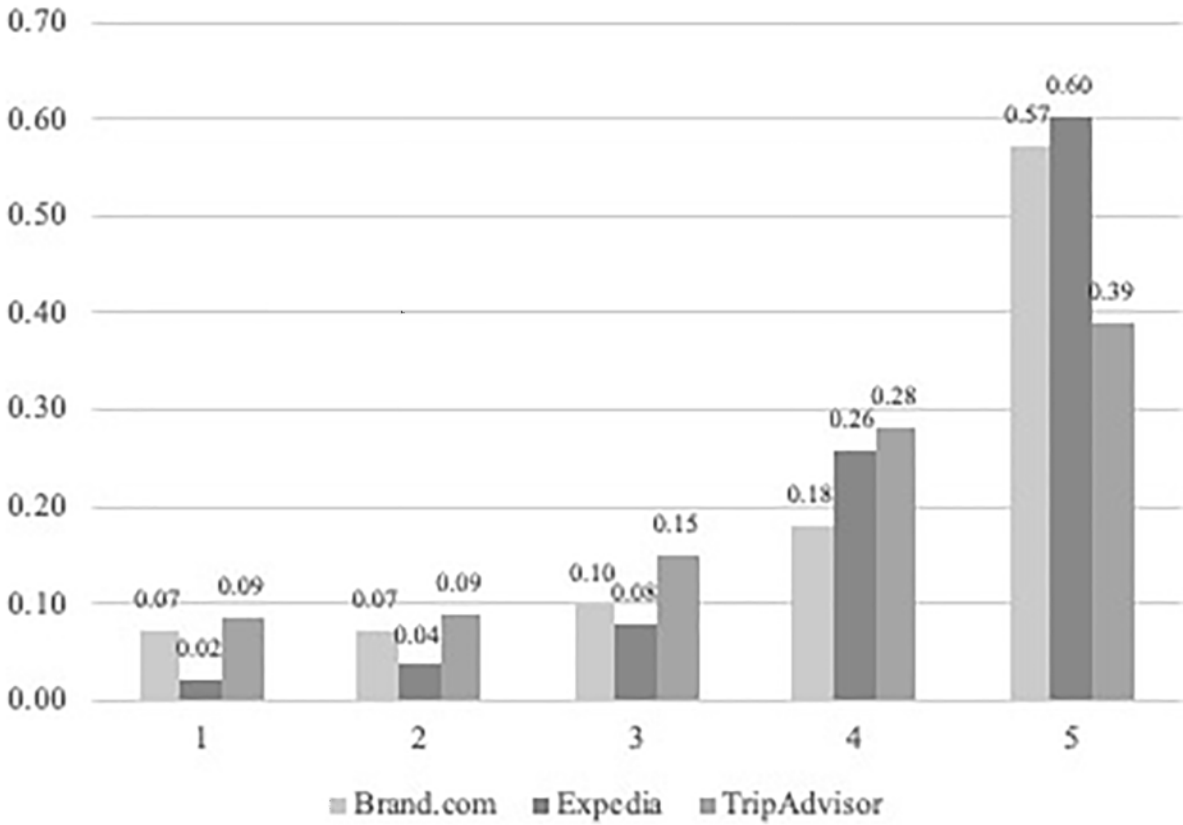
**Table 4.**  
Top 10 Most Frequently Used Words for the Same Hotel's Reviews by Platforms.

	TripAdvisor	Expedia	Hotel Brand
1	room (1,346)	room (754)	location (678)
2	time (1,035)	location (691)	room (654)
3	stay (659)	time (482)	stay (644)
4	square (536)	staff (400)	time (562)
5	location (463)	square (334)	staff (488)
6	get (367)	stay (288)	square (352)
7	service (360)	leave (287)	service (295)
8	one (349)	review (282)	excellent (228)
9	staff (336)	comments (281)	good (217)
10	floor (331)	traveler (281)	clean (210)

reviews may be more extreme (and perhaps nonrepresentative; Anderson et al., 1994). In contrast, reviews that were more easily posted (due to the email invitation) are less likely to be affected by the different posting motivation across different ratings. As a result, reviewers that would otherwise have relatively low posting intention (i.e., positive ratings) are able to easily submit their reviews, potentially resulting in higher average ratings at sites with email encouraged/prompted reviews, for example, OTAs and brand sites (Han & Anderson, 2020).

### Price Differences Across Platforms

Next to reviews, prices are among the most interesting data types available across different travel platforms. Just as



**Figure 7.**  
Distribution of Review Scores by Platform.



**Figure 8.**  
Listing of a Sold-Out Property in Expedia.

with reviews, different platforms have different objectives and different capabilities when presenting prices. As a result, the way in which prices are displayed and communicated across platforms varies, which may influence customers' purchasing decisions. For example, Expedia displays all properties, including those that are sold-out or have limited availability, as is shown in Figure 8. In contrast, other travel platforms, such as Airbnb, only list properties that are available and have price data to be displayed. Customers may behave differently as a result of these differences in

information display. Accordingly, it is likely that customers behave differently depending on which travel platform they choose (Park & Jang, 2018). Therefore, assuming that all platforms have the same market structure ignores these differences in information stimuli that influence customer behavior.

OTAs are also considerably more versed in merchandising when displaying price information. Merchandising by OTAs may include such actions as strike-through pricing or scarcity messaging, which attempt to increase conversion

Member Rate Advance Purchase, prepay in full, non-refundable if cancelled more than 1 day after booking, no changes, see Rate details

[Rate Details](#)

259  
**184** USD / night

**SELECT**



Get a \$250 statement credit.

THIS STAY COULD BE \$-66.

[See Details.](#) ✓

Advance Purchase rate, prepay in full, non-refundable if cancelled less than 1 day before arrival, no changes, see Rate details

[Rate Details](#)

259  
**194** USD / night

**SELECT**

**Figure 9.**  
**Different Price Offer to the Members.**

rates. Similarly, rate parity has received considerable attention recently. While OTAs want to display prices that can compete with those posted on supplier websites, hotels are often motivated to offer better prices to customers who book directly to avoid OTA commissions. Closed user group (CUG) or membership selling is a common method of offering lower prices to a subset of consumers. In CUG situations, the consumer must login (either at the OTA or hotel site) to access better prices. Figure 9 illustrates strike-through prices and the resulting member rates available to a typical CUG. By definition, CUG discounted prices are only available to those who have logged in; the majority of customers will never see this price disparity. Researchers who analyze the online hotel market using prices scraped from travel websites must understand how these distribution channels influence each other.

## Implications

The academic implications of web scraping are twofold: exploratory and confirmatory. In general, web scraped data are suitable for exploratory research where the research question has not been examined in detail. The aim of this research area is to understand the general pattern of the customers in question and find preliminary evidence that warrants a more detailed study. For example, Danescu-Niculescu-Mizil et al. (2013) show that users in online beer communities follow a two-stage life cycle in terms of the language they use: the innovative learning phase and the conservative phase. Wu and Huberman (2010) found that later online product review ratings tend to vary considerably from earlier ones, making overall review ratings less extreme. Although these studies demonstrate interesting behavioral patterns and

are valuable in their own right, they still require further testing in an experimental setting.

Web scraped data are also often used in confirmatory studies that test hypotheses from previous studies (Landers et al., 2016; Marres & Weltevrede, 2013). If this hypothesis testing investigates whether one of the variables scraped from the website is an exogenous variable that varies independently of an error term, this study is considered as an experiment (Harrison & List, 2004). Web scraped data enable researchers to conduct experiments in a real environmental context and conduct research without informing research subjects that they are taking part in an experiment. This research design is referred to as a natural experiment. The advantage of a natural experiment is that it accounts for the realistic environment that real customers encounter. Therefore, an ideal natural experiment not only increases external validity but also does so when internal validity is insufficient (Harrison & List, 2004). For example, Han and Anderson (2020) take advantage of a unique characteristic of TripAdvisor, the fact that a portion of its reviewers post after being prompted to do so by hotel managers. By comparing regular self-motivated reviews and prompted reviews, the authors test whether satisfied or dissatisfied customers are more motivated to post online reviews. Web scraped data can also be combined with data from other sources in confirmatory research. Xie et al. (2014) combined TripAdvisor's review data with archival data regarding hotel revenue per available room (RevPAR) matched to the Texas Comptroller's Office database to investigate the impact of various review website attributes on hotel performance.

Managers in the hospitality industry can benefit from web scraping, too. It is well known that customers heavily rely on online reviews before making a purchasing decision



(Brown et al., 2007; Chevalier & Mayzlin, 2006). Therefore, it is important to understand how customers evaluate services. Although there are multiple online reputation management companies that do this job on behalf of firms, hotels can save money by web scraping themselves. For example, firms can web scrape major online review websites in the industry periodically and automate this process, which is referred to as web crawling (Massimino, 2016). Without spending any extra money, firms can obtain a summary of each online review website. Based on these summaries, managers can decide which review websites they should prioritize and invest to maximize the number of reviews they receive.

## Ethical Concerns

As legislation on web scraping varies from country to country, researchers should look into local legislation. In this section, we discuss the legal and ethical issues mainly in the U.S. context. On September 9, 2019, the U.S. Supreme Court legalized web scraping in situations where the scraped information is designed to be publicly accessible. The court defined public information as data that are neither available for purchase nor hidden behind a password-protected authentication system. The logic behind the court's decision was that, legally, web scraping is no different than browsing in terms of what data are being requested from a website. However, web scraping information that is accessible exclusively to the members and requires logging in is illegal, as this behavior explicitly violates the terms of service (ToS).

It is also noteworthy that web scraping copyrighted data and re-using them for commercial purposes would be considered illegal. For example, web scraping video contents from YouTube and re-posting them on ones' own website could be illegal as videos are copyrighted.

However, illegally sharing data is not likely a matter of concern for the majority of our target audience: researchers who are interested in using web scraping for academic purposes. In addition, web scraping is a relatively new data collection method in academia, and therefore the law is still evolving (Hillen, 2019). Therefore, researchers must take into account that it is always possible that current laws regarding web scraping will change and that they may need to seek professional legal advice before web scraping.

In addition to legal issues, researchers should also consider ethical constraints when collecting data online for academic purposes. As previous studies argue, legality does not necessarily mean that data usage is entirely ethical (Massimino, 2016). Using online data in research is relatively new in comparison to other data sources. Therefore, its current legality suggests a need for further research regarding the ethicality of web scraped data and the safety

of the web scraping practice. By the same token, although web scraping qualifies for Institutional Review Boards (IRBs)' review exemptions in most of the cases (Massimino, 2016), researchers need to be conscientious of any societal entities that may be impacted by web scraping. A common ethical concern regarding web scraping is related to the problem of sending too many requests to the host over a short span of time. A typical web scraper involves querying a website repeatedly. If overused, this practice can prevent others from accessing the website. A web scraper that is written for the purpose of collecting multiple online reviews or hotel prices sends requests to the web server that is hosting the site whenever it opens a new page. While the requests of a human user are usually within a manageable range, a web scraper that makes speedy and bulky automated requests can easily exceed the bandwidth threshold of the host and make the server unresponsive (Massimino, 2016). When the web scraper hits the server with frequent requests, a host may issue a warning or may respond with useless content if web scraping behavior is detected.

Therefore, we recommend inserting a random delay between individual requests, such as limiting requests to three per second (Massimino, 2016). For example, Landers et al. (2016) executed a 2-s delay between each web page request to avoid overburdening the host's server. In addition, scraping during off-peak hours can help reduce the load on the host and increase the speed of the scraping process. Finally, a good practice in web scraping is to carefully read the robots exclusion protocol (REP), which are standardized instructions on whether certain user agents can scrape parts of a website. Typically, this information can be found in the admin page of a website.

## Summary

A growing number of customers not only obtain travel information online but also make transactions over the internet. Studies indicate that the internet (as opposed to the offline, voice, or travel agent distribution channels) has become the dominant distribution channel in terms of travel reservations (Park, 2009). Therefore, the importance of studying online customer behaviors cannot be overemphasized. However, there is a strong tendency among hospitality researchers to rely on traditional data collection methods, which are limited in terms of what research questions they can answer as well as the generalizability of their insights. Our work provides a simple method of scraping online reviews and price data using the Python language. As our goal is to make hospitality researchers more comfortable collecting online data, we focus on how the scraping process functions on major travel websites. Although not a comprehensive introduction into Python, we introduce the essential elements necessary to handle and scrape interactive

hospitality platforms such that they can augment readily available introductory Python resources.

Although web scraped data introduce incredible opportunities to hospitality researchers, there are important aspects that must be accounted for when scraping travel websites. As every platform has unique characteristics and purposes, each platform attracts different users, which in turn forms a different market. While ignoring these platform-specific aspects may induce biases in the researchers' analyses, properly making use of these challenges may make platform differences a unique natural experiment setting for exploring new opportunities. At the same time, embracing these differences is what provides for fruitful research. For instance, the social influence effect, which refers to the effect of previous reviews on future reviews, impacts most online review platforms due to the nature of being able to see previous reviewers' opinions. This effect induces potential biases that are not of concern in traditional survey research as there is no chance that survey participants will see other participants' opinions before answering the survey questions. Despite the possibility of inducing social influence bias, Askalidis et al. (2017) overcome this challenge by utilizing reviews written by retailer-prompted reviewers who were invited to contribute their opinions using a separate web page where there is a lesser chance of seeing previous opinions. Going further, they compare this group of reviews to the regular, organic reviews using the difference-in-difference method, and make use of the challenge to identify the social influence effect in the online communities. Another example of turning the challenge into opportunity is the study of Wang and Chaudhry (2018). Many online hotel review platforms allow managers to respond to customer reviews, which may influence future reviewers' opinions. While the managerial response could present a challenge for researchers who want to understand the unbiased opinions of the reviewers, Wang and Chaudhry (2018) identified the managerial response effect by comparing ratings from online review platforms where managerial responses are visible with ratings from platforms where managerial responses are not made visible.

Our article has significant implications for hospitality researchers who hope to better understand the online travel marketplace. We outline simple methods that enable hospitality researchers to collect incredibly useful secondary data that they could not have obtained by relying on traditional data collection methods alone. Although traditional data collection methods are still valuable and, in many cases, cannot be replaced by new methods, online customer behaviors are hardly replicable in offline research design. Even if the purpose of the research is to make a causal inference that can only be tested in a strict lab setting, confirming this effect in the real online marketplace adds value to the study in terms of external validity.

## Declaration of Conflicting Interests

The author(s) declared no potential conflicts of interest with respect to the research, authorship, or publication of this article.

## Funding

The author(s) received no financial support for the research, authorship, or publication of this article.

## ORCID iDs

Saram Han  <https://orcid.org/0000-0001-9077-9107>

Christopher K. Anderson  <https://orcid.org/0000-0002-7103-8351>

## Supplemental Material

Supplemental material for this article is available online.

## References

- Anderson, E. W., Fornell, C., & Lehmann, D. R. (1994). Customer satisfaction, market share, and profitability: Findings from Sweden. *Journal of Marketing*, 58(3), 53–66.
- Askalidis, G., Kim, S. J., & Malthouse, E. C. (2017). Understanding and overcoming biases in online review systems. *Decision Support Systems*, 97, 23–30.
- Brown, J., Broderick, A. J., & Lee, N. (2007). Word of mouth communication within online communities: Conceptualizing the online social network. *Journal of Interactive Marketing*, 21(3), 2–20. <https://doi.org/10.1002/dir.20082>
- Chevalier, J. A., Dover, Y., & Mayzlin, D. (2018). Channels of impact: User reviews when quality is dynamic and managers respond. *Marketing Science*, 37, 685–853.
- Chevalier, J. A., & Mayzlin, D. (2006). The effect of word of mouth on sales: Online book reviews. *Journal of Marketing Research*, 43(3), 345–354. <https://doi.org/10.1509/jmkr.43.3.345>
- Choi, S., Mattila, A. S., Van Hoof, H. B., & Quadri-Felitti, D. (2017). The role of power and incentives in inducing fake reviews in the tourism industry. *Journal of Travel Research*, 56(8), 975–987.
- Danescu-Niculescu-Mizil, C., West, R., Jurafsky, D., Leskovec, J., & Potts, C. (2013). No country for old members. *Proceedings of the 22nd International Conference on World Wide Web - WWW '13, Rio de Janeiro, Brazil* (pp. 307–318). Association for Computing Machinery. <https://doi.org/10.1145/2488388.2488416>
- Downey, A. (2014). *Think Python: How to think like a computer scientist*. Green Tea Press.
- Ert, E., & Fleischer, A. (2019). The evolution of trust in Airbnb: A case of home rental. *Annals of Tourism Research*, 75, 279–287.
- Gao, B., Hu, N., & Bose, I. (2017). Follow the herd or be myself? An analysis of consistency in behavior of reviewers and helpfulness of their reviews. *Decision Support Systems*, 95, 1–11.
- Han, S., & Anderson, C. K. (2020). Customer motivation and response bias in online reviews. *Cornell Hospitality*

- Quarterly*, 61, 142–153. <https://doi.org/10.1177/1938965520902012>
- Harrison, G. W., & List, J. A. (2004). Field experiments. *Journal of Economic Literature*, 42(4), 1009–1055.
- Hillen, J. (2019). Web scraping for food price research. *British Food Journal*, 121(12), 3350–3361. <https://doi.org/10.1108/BFJ-02-2019-0081>
- Kupor, D., & Tormala, Z. (2018). When moderation fosters persuasion: The persuasive power of deviatory reviews. *Journal of Consumer Research*, 45, 490–510.
- Landers, R. N., Brusso, R. C., Cavanaugh, K. J., & Collmus, A. B. (2016). A primer on theory-driven web scraping: Automatic extraction of big data from the internet for use in psychological research. *Psychological Methods*, 21(4), 475–492.
- Litvin, S. W., & Sobel, R. N. (2019). Organic versus solicited hotel TripAdvisor reviews: Measuring their respective characteristics. *Cornell Hospitality Quarterly*, 60(4), 370–377.
- Marres, N., & Weltevrede, E. (2013). Scraping the social?: Issues in live social research. *Journal of Cultural Economy*, 6(3), 313–335.
- Massimino, B. (2016). Accessing online data: Web-crawling and information-scraping techniques to automate the assembly of research data. *Journal of Business Logistics*, 37(1), 34–42.
- Mayzlin, D., Dover, Y., & Chevalier, J. (2014). Promotional reviews: An empirical investigation of online review manipulation. *American Economic Review*, 104(8), 2421–2455.
- Min, H., Lim, Y., & Magnini, V. P. (2015). Factors affecting customer satisfaction in responses to negative online hotel reviews: The impact of empathy, paraphrasing, and speed. *Cornell Hospitality Quarterly*, 56(2), 223–231.
- Park, J. (2009). Consumers' travel website transferring behaviour: Analysis using clickstream data-time, frequency, and spending. *Service Industries Journal*, 29(10), 1451–1463.
- Park, J. Y., & Jang, S. C. S. (2018). The impact of sold-out information on tourist choice decisions. *Journal of Travel and Tourism Marketing*, 35(5), 622–632.
- Schahn, J., & Holzer, E. (1990). Studies of individual environmental concern. *Environment and Behavior*, 22(6), 767–786.
- Shin, H., Perdue, R. R., & Pandelaere, M. (2019). Managing customer reviews for value co-creation: An empowerment theory perspective. *Journal of Travel Research*, 59, 792–810. <https://doi.org/10.1177/0047287519867138>
- vanden Broucke, S., & Baesens, B. (2018). *Practical web scraping for data science*. Apress.
- Viglia, G., & Dolnicar, S. (2020). A review of experiments in tourism and hospitality. *Annals of Tourism Research*, 80, 102858.
- Wang, Y., & Chaudhry, A. (2018). When and how managers' responses to online reviews affect subsequent reviews. *Journal of Marketing Research*, 55(2), 163–177.
- Wu, F., & Huberman, B. A. (2010). Opinion formation under costly expression. *ACM Transactions on Intelligent Systems and Technology*, 1(1), 1–13. <https://doi.org/10.1145/1858948.1858953>
- Wu, L., Mattila, A. S., Wang, C. Y., & Hanks, L. (2015). The impact of power on service customers' willingness to post online reviews. *Journal of Service Research*, 19(2), 224–238.
- Xie, K. L., Zhang, Z., & Zhang, Z. (2014). The business value of online consumer reviews and management response to hotel performance. *International Journal of Hospitality Management*, 43, 1–12. <https://doi.org/10.1016/j.ijhm.2014.07.007>
- Zervas, G., Proserpio, D., & Byers, J. (2018). A first look at online reputation on Airbnb, where every stay is above average. *SSRN Electronic Journal*. <http://dx.doi.org/10.2139/ssrn.2554500>
- Zhang, Y., Liu, L., & Ho, S. Y. (2020). How do interruptions affect user contributions on social commerce? *Information Systems Journal*, 30(3), 535–565.

### Author Biographies

**Saram Han** is an assistant professor at the College of Business and Technology at Seoul National University of Science and Technology. He received his PhD from the Cornell School of Hotel Administration in the Cornell SC Johnson College of Business. His research interests include digital marketing, online reviews, service marketing, and marketing analytics. He earned a BBA in Tourism Management from Kyung-hee University, Seoul, Korea, and an M.S. degree from the Michigan Program in Survey Methodology, University of Michigan.

**Christopher K. Anderson** is a professor at Cornell SC Johnson College of Business, School of Hotel Administration. He earned his BSc/MSc in engineering from the University of Guelph, and his MBA/PhD from the University of Western Ontario, Richard Ivey School of Business. He teaches and conducts research in data analytics, pricing, distribution, and revenue management.